



Consistent change-point detection with kernels

Damien Garreau, Sylvain Arlot

► To cite this version:

Damien Garreau, Sylvain Arlot. Consistent change-point detection with kernels. *Electronic Journal of Statistics* , 2018, 12 (2), pp.4440-4486. hal-01416704v2

HAL Id: hal-01416704

<https://hal.science/hal-01416704v2>

Submitted on 28 Jun 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Consistent change-point detection with kernels

Damien Garreau*

*Centre de recherche Inria de Paris
2 rue Simone Iff
CS 42112
75589 Paris Cedex 12
e-mail: damien.garreau@ens.fr*

and

Sylvain Arlot

*Laboratoire de Mathématiques d'Orsay
Univ. Paris-Sud, CNRS, Université Paris-Saclay,
91405 Orsay, France
e-mail: sylvain.arlot@math.u-psud.fr*

Abstract: In this paper we study the kernel change-point algorithm (KCP) proposed by Arlot, Celisse and Harchaoui [4], which aims at locating an unknown number of change-points in the distribution of a sequence of independent data taking values in an arbitrary set. The change-points are selected by model selection with a penalized kernel empirical criterion. We provide a non-asymptotic result showing that, with high probability, the KCP procedure retrieves the correct number of change-points, provided that the constant in the penalty is well-chosen; in addition, KCP estimates the change-points location at the optimal rate. As a consequence, when using a characteristic kernel, KCP detects all kinds of change in the distribution (not only changes in the mean or the variance), and it is able to do so for complex structured data (not necessarily in \mathbb{R}^d). Most of the analysis is conducted assuming that the kernel is bounded; part of the results can be extended when we only assume a finite second-order moment.

MSC 2010 subject classifications: Primary 62M10; secondary 62G20.

Keywords and phrases: change-point detection, kernel methods, penalized least-squares.

Received July 2017.

1. Introduction

In many situations, some properties of a time series change over time, such as the mean, the variance or higher-order moments. Change-point detection is the long standing question of finding both the number and the localization of such changes. This is an important front-end task in many applications. For instance, detecting changes occurring in comparative genomic hybridization array

*Corresponding author

data (CGH arrays) is crucial to the early diagnosis of cancer [34]. In finance, some intensively examined time series like the volatility process exhibit local homogeneity and it is useful to be able to segment these time series both for modeling and forecasting [38, 49]. Change-point detection can also be used to detect changes in the activity of a cell [45], in the structure of random Markov fields [43], or a sequence of images [30, 1]. Generally speaking, it is of interest to the practitioner to segment a time series in order to calibrate its model on homogeneous sets of datapoints.

Addressing the change-point problem in practice requires to face several important challenges. First, the number of changes can not be assumed to be known in advance — in particular, it can not be assumed to be equal to 0 or 1 —, hence a practical change-point procedure must be able to infer the number of changes from the data. Second, changes do not always occur in the mean or the variance of the data, as assumed by most change-point procedures. We need to be able to detect changes in other features of the distribution. Third, parametric assumptions — which are often made for building or for analyzing change-point procedures — are often unrealistic, so that we need a fully non-parametric approach. Fourth, data points in the time series we want to segment can be high-dimensional and/or structured. If the dimensionality is larger than the number of observations, a non-asymptotic analysis is mandatory for theoretical results to be meaningful. When data are structured — for instance, histograms, graphs or strings —, taking their structure into account seems necessary for detecting efficiently the change-points.

We focus only on the *offline* problem in this article, that is, when all observations are given at once, as opposed to the situation where data come as a continuous stream. We refer to Tartakovsky, Nikiforov and Basseville [51] for an extensive review of sequential methods, which are adapted to the latter situation. Numerous offline change-point procedures have been proposed since the seminal works of Page [44], Fisher [21] and Bellman [9], which are mostly parametric in essence. We refer to Brodsky and Darkhovsky [13, Chapter 2] for a review of non-parametric offline change-point detection methods. Among recent works in this direction, we can mention the Wild Binary Segmentation (WBS, [22]) and the non-parametric multiple change-point detection procedure (NMCD, [56]). Some authors also consider the case of high-dimensional data when only a few coordinates of the mean change at each change-point [53, and references therein], or the problem of detecting gradual changes [52]; this paper does not address these slightly different problems.

To the best of our knowledge, no offline change-point procedure addressed simultaneously the four challenges mentioned above, until the kernel change-point procedure (KCP) was proposed by Arlot, Celisse and Harchaoui [4]. In short, KCP mixes the penalized least-squares approach to change-point detection [17, 39] with semi-definite positive kernels [5]. It is not the only procedure that uses positive semi-definite kernels to detect changes in a times series. Apart from Harchaoui and Cappé [27], who introduced KCP for a fixed number of change-point, and Arlot, Celisse and Harchaoui [4] who extended KCP to an unknown number of change-points, we are aware of several closely related work.

Maximum Mean Discrepancy [MMD, 24] has been used for building two sample tests; a block average version of the MMD, named the M -statistic, has lead to an online change-point detection procedure [41]. A kernel-based statistic, named kernel Fisher discriminant ratio, has been used by Harchaoui, Moulines and Bach [28] for homogeneity testing and for detecting one change-point. Sharipov, Tewes and Wendler [48] build an analogue of the CUSUM statistic for Hilbert-valued random variables in order to detect a single change in the mean, and could be applied in our setting to the images of the observations in the feature space. Kernel change detection [18] is an online procedure that uses a kernel to build a dissimilarity measure between the near past and future of a data-point.

On the computational side, the KCP segmentation can be computed efficiently thanks to a dynamic programming algorithm [27, 4], which can be made even faster [16]. An oracle inequality for KCP is proved by Arlot, Celisse and Harchaoui [4]; this is not exactly a result on change-point estimation, but a guarantee on estimation of the “mean” of the time series in the RKHS associated with the kernel chosen. The good numerical performance of KCP in terms of change-point estimation is also demonstrated in several experiments.

So, a key theoretical question remains open: does KCP estimate correctly the number of change-points and their locations with a large probability? If yes, at which speed does KCP estimate the change-point locations?

This paper answers these questions, showing that KCP has good theoretical properties for change-point estimation with independent data, under a boundedness assumption (Theorem 3.1 in Section 3.1). This result is non-asymptotic, hence meaningful for high-dimensional or complex data. In the asymptotic setting — with a fixed true segmentation and more and more data points observed within each segment —, Theorem 3.1 implies that KCP estimates consistently all changes in the “kernel mean” of the distribution of data, at speed $\log(n)/n$ with respect to the sample size n . Since we make no assumptions on the minimal size of the true segments, this matches minimax lower bounds [14]. We also provide a partial result under a weaker finite variance assumption (Theorem 3.2 in Section 3.3) and explain in Section 5 how our proofs could be extended to other settings, including the dependent case. These findings are illustrated by numerical simulations in Section 4.

An important case is when KCP is used with a characteristic kernel [23], such as the Gaussian or the Laplace kernel. Then, any change in the distribution of data induces a change in the “kernel mean”. So, Theorem 3.1 implies that KCP then estimates consistently and at the minimax rate *all changes* in the distribution of the data, without any parametric assumption and without prior knowledge about the number of changes.

Our results also are interesting regarding to the theoretical understanding of least-squares change-point procedures. Indeed, when KCP is used with the linear kernel, it reduces to previously known penalized least-squares change-point procedures [54, 17, 39, for instance]. There are basically two kinds of results on such procedures in the change-point literature: (i) asymptotic statements on change-point estimation [54, 55, 6, 37] and (ii) non-asymptotic oracle inequalities

[17, 39, 4], which are based upon concentration inequalities and model selection theory [10] but do not directly provide guarantees on the estimated change-point locations. Our results and their proofs show how to reconcile the two approaches when we are interested in change-point locations, which is already new for the case of the linear kernel, and also holds for a general kernel.

2. Kernel change-point detection

This section describes the general change-point problem and the kernel change-point procedure [4].

2.1. Change-point problem

Set $2 \leq n < +\infty$ and consider X_1, \dots, X_n independent \mathcal{X} -valued random variables, where \mathcal{X} is an arbitrary (measurable) space. The goal of change-point detection is to detect abrupt changes in the distribution of the X_i s. For any $D \in \{1, \dots, n\}$ and any integers $0 = \tau_0 < \tau_1 < \dots < \tau_D = n$, we define the *segmentation* $\tau := [\tau_0, \dots, \tau_D]$ of $\{1, \dots, n\}$ as the collection of segments $\lambda_\ell = \{\tau_{\ell-1} + 1, \dots, \tau_\ell\}$, $\ell \in \{1, \dots, D\}$. We call *change-points* the right-end of the segments, that is the τ_ℓ , $\ell \in \{1, \dots, D\}$. We denote by \mathcal{T}_n^D the set of segmentations with D segments and $\mathcal{T}_n := \bigcup_{D=1}^n \mathcal{T}_n^D$ the set of all segmentations of $\{1, \dots, n\}$. For any $\tau \in \mathcal{T}_n$, we write D_τ for the number of segments of τ . Figure 1 provides a visual example.



FIG 1. We often represent the segmentations as above. The bullet points stand for the elements of $\{1, \dots, n\}$. Here, $n = 10$, $D_\tau = 3$, $\tau_0 = 0$, $\tau_1 = 3$, $\tau_2 = 7$ and $\tau_3 = 10$.

An important example to have in mind is the following.

Example 2.1 (Asymptotic setting). Let $K \geq 1$, $0 = b_0 < b_1 < \dots < b_K < b_{K+1} = 1$ and P_1, \dots, P_{K+1} some probability distributions on \mathcal{X} be fixed. Then, for any n and $i \in \{1, \dots, n\}$, we set $t_i := i/n$ and the distribution of X_i is $P_{j(i)}$ where $j(i)$ is such that $t_i \in [b_j, b_{j+1})$. In other words, we have a fixed segmentation of $[0, 1]$, given by the b_j , a fixed distribution over each segment, given by the P_j , and we observe independent realizations from the distributions at discrete times t_1, \dots, t_n . The corresponding true change-points in $\{0, \dots, n\}$ are the $\lfloor nb_j \rfloor$, $j = 1, \dots, K$. For n large enough, there are $K + 1$ segments. Figure 2 shows an example. Let us emphasize that in this setting, n going to infinity does not mean that new observations are observed over time. Recall that we consider the change-point problem *a posteriori*: a larger n means that we have been able to observe the phenomenon of interest with a finer time discretization. Also note that this asymptotic setting is restrictive in the sense that segments size asymptotically are of order n ; we do not make this assumption

in our analysis, which also covers asymptotic settings where some segments have a smaller size.

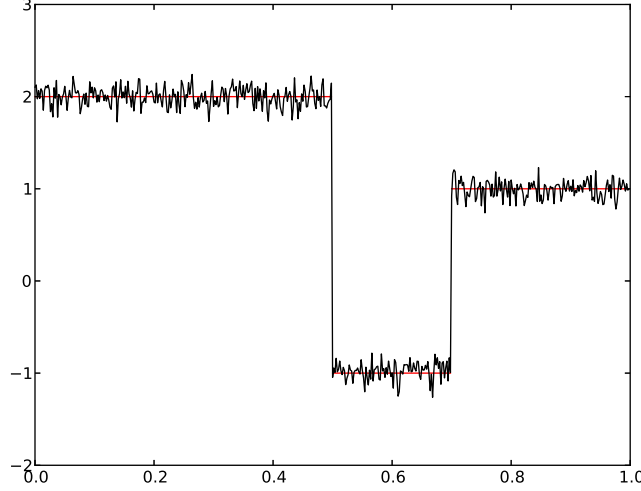


FIG 2. Illustration of the asymptotic setting (Example 2.1) in the case of changes in the mean of the X_i . Here, $\mathcal{X} = \mathbb{R}$, $X_i = f(t_i) + \varepsilon_i$ with $\varepsilon_1, \dots, \varepsilon_n$ i.i.d. and centered, and $f : [0, 1] \rightarrow \mathbb{R}$ is a (fixed) piecewise constant function (shown in red). The goal is to recover the number of abrupt changes of f (here, 2) and their locations ($b_1 = 0.5$ and $b_2 = 0.7$). Note that other kinds of changes in the distribution of the X_i can be considered, see Section 4.

2.2. Kernel change-point procedure (KCP)

Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a positive semidefinite kernel, that is, a measurable function such that the matrix $(k(x_i, x_j))_{1 \leq i, j \leq m}$ is positive semidefinite for any $m \geq 1$ and $x_1, \dots, x_m \in \mathcal{X}$ [46]. Classical examples of kernels are given by [4, section 3.2], among which:

- the *linear kernel*: $k^{\text{lin}}(x, y) = \langle x, y \rangle_{\mathbb{R}^p}$ for $x, y \in \mathcal{X} = \mathbb{R}^p$.
- the *polynomial kernel* of order $d \geq 1$: $k_d^{\text{poly}}(x, y) = (\langle x, y \rangle_{\mathbb{R}^p} + 1)^d$ for $x, y \in \mathcal{X} = \mathbb{R}^p$.
- the *Gaussian kernel* with bandwidth $h > 0$: $k_h^{\text{G}}(x, y) = \exp[-\|x - y\|^2 / (2h^2)]$ for $x, y \in \mathcal{X} = \mathbb{R}^p$.
- the *Laplace kernel* with bandwidth $h > 0$: $k_h^{\text{L}}(x, y) = \exp[-\|x - y\| / (2h^2)]$ for $x, y \in \mathcal{X} = \mathbb{R}^p$.
- the χ^2 -kernel: $k_{\chi^2}(x, y) = \exp\left(-\frac{1}{2} \sum_{i=1}^p \frac{(x_i - y_i)^2}{x_i + y_i}\right)$ for $x, y \in \mathcal{X}$ the p -dimensional simplex.

As done by Harchaoui and Cappé [27] and Arlot, Celisse and Harchaoui [4], for a given segmentation $\tau \in \mathcal{T}_n^D$, we assess the adequation of τ with the *kernel least-squares criterion*

$$\begin{aligned} \widehat{\mathcal{R}}_n(\tau) &:= \frac{1}{n} \sum_{i=1}^n k(X_i, X_i) \\ &- \frac{1}{n} \sum_{\ell=1}^D \left[\frac{1}{\tau_\ell - \tau_{\ell-1}} \sum_{i=\tau_{\ell-1}+1}^{\tau_\ell} \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} k(X_i, X_j) \right]. \end{aligned} \quad (2.1)$$

Elementary algebra shows that, when $\mathcal{X} = \mathbb{R}^p$ and $k = k^{\text{lin}}$, $\widehat{\mathcal{R}}_n$ is the usual least-squares criterion. Minimizing this criterion over the set of all segmentations always outputs the segmentation with n segments reduced to a point, that is $[0, \dots, n]$; this is a well-known overfitting phenomenon. To counteract this, a classical idea [36, for instance] is to minimize a penalized criterion $\text{crit}(\tau) := \widehat{\mathcal{R}}_n(\tau) + \text{pen}(\tau)$, where $\text{pen} : \mathcal{T}_n \rightarrow \mathbb{R}_+$ is called the penalty. Formally, the kernel change-point procedure (KCP) of Arlot, Celisse and Harchaoui [4] selects the segmentation

$$\widehat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n} \{ \text{crit}(\tau) \} \quad \text{where} \quad \text{crit}(\tau) = \widehat{\mathcal{R}}_n(\tau) + \text{pen}(\tau). \quad (2.2)$$

In this paper, we focus on the classical choice of a penalty proportional to the number of segments, similarly to AIC, BIC and C_p criteria. Namely, we consider

$$\text{pen}(\tau) = \text{pen}_\ell(\tau) := \frac{CM^2 D_\tau}{n}, \quad (2.3)$$

where C is a positive constant and M is specified in Assumption 1 later on. As mentioned in the Introduction, slightly different penalty shapes can be considered, as suggested by Arlot, Celisse and Harchaoui [4]. Our results could be extended to the penalty of Arlot, Celisse and Harchaoui [4], but we choose to consider the linear penalty (2.3) only for simplicity.

2.3. The reproducing kernel Hilbert space

Let \mathcal{H} be the reproducing kernel Hilbert space (RKHS) associated with k [5], together with the canonical feature map $\Phi : \mathcal{X} \rightarrow \mathcal{H}$

$$\begin{aligned} \Phi &: \mathcal{X} \rightarrow \mathcal{H} \\ x &\mapsto \Phi(x) := k(\cdot, x). \end{aligned}$$

We write $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ (resp. $\|\cdot\|_{\mathcal{H}}$) for the inner product (resp. the norm) of \mathcal{H} . For any $i \in \{1, \dots, n\}$, define $Y_i := \Phi(X_i) \in \mathcal{H}$. In the case where $k = k^{\text{lin}}$, then $Y_i = \langle \cdot, X_i \rangle_{\mathbb{R}^p}$ and the empirical risk $\widehat{\mathcal{R}}_n$ reduces to the least-squares criterion

$$\widehat{\mathcal{R}}_n(\tau) = \frac{1}{n} \sum_{\ell=1}^{D_\tau} \sum_{i=\tau_{\ell-1}+1}^{\tau_\ell} (X_i - \bar{X}_\ell)^2,$$

where \bar{X}_ℓ is the empirical mean of the X_i over the segment $\{\tau_{\ell-1} + 1, \dots, \tau_\ell\}$. It is well-known that penalized least-squares procedures detect changes in the mean of the observations X_i , see Yao [54]. Hence the kernelized version of this least-squares procedure, KCP, should detect changes in the “mean” of the $Y_i = \Phi(X_i)$, which are a nonlinear transformation of the X_i .

More precisely, assume that \mathcal{H} is separable and that

$$\forall i \in \{1, \dots, n\}, \quad \mathbb{E} \left[\sqrt{k(X_i, X_i)} \right] < +\infty.$$

Then μ_i^* , the Bochner integral of Y_i , is well-defined [40]. The condition above is satisfied in our setting (when either Assumption 1 or Assumption 2 holds true, see Section 2.5), and \mathcal{H} is separable in most cases [20]. The Bochner integral commutes with continuous linear operators, hence the following property holds, which will be of common use:

$$\forall g \in \mathcal{H}, \quad \langle \mu_i^*, g \rangle_{\mathcal{H}} = \mathbb{E}[g(X_i)] = \mathbb{E}[\langle Y_i, g \rangle_{\mathcal{H}}].$$

We now define the “true segmentation” $\tau^* \in \mathcal{T}_n$ by

$$\begin{aligned} \mu_1^* = \dots = \mu_{\tau_1^*}^*, \quad \mu_{\tau_1^*+1}^* = \dots = \mu_{\tau_2^*}^*, \quad \dots \quad \mu_{\tau_{D_{\tau^*}-1}^*+1}^* = \dots = \mu_n^* \\ \text{and } \forall i \in \{1, \dots, D_{\tau^*} - 1\}, \quad \mu_{\tau_i^*}^* \neq \mu_{\tau_{i+1}^*}^* \end{aligned} \quad (2.4)$$

with $1 \leq \tau_1^* < \dots < \tau_{D_{\tau^*}-1}^* \leq n$. We call the τ_i^* s the *true* change-points. It should be clear that it is always possible to define τ^* .

A kernel is said to be characteristic if the mapping $P \mapsto \mathbb{E}_{X \sim P} [\Phi(X)]$ is injective, for P belonging to the set of Borel probability measures on \mathcal{X} [50]. In simpler terms, when k is a characteristic kernel, X_i and X_{i+1} have the same distribution if and only if $\mu_i^* = \mu_{i+1}^*$, and τ^* indeed corresponds to the set of changes in the distribution of the X_i . For instance, all strictly positive definite kernels are characteristic, including the Gaussian kernel, see Sriperumbudur et al. [50]. Therefore, in the setting of Example 2.1, for n large enough, $D_{\tau^*} = K + 1$ and $\tau_\ell^* = \lfloor nb_\ell \rfloor$ for $\ell = 1, \dots, K$.

For a general kernel, some changes of P_{X_i} , the distribution of X_i , might not appear in τ^* . For instance, with the linear kernel, τ^* only corresponds to changes of the mean of the X_i . In most cases, a characteristic kernel is known and we can choose to use KCP with a characteristic kernel; then, as we prove in the following, KCP eventually detects any change in the distribution of the observations. But one can also choose a non-characteristic kernel on purpose, hence focusing only on some changes in the distribution of the X_i . For instance, the polynomial kernel of order d is not characteristic and leads to the detection of changes in the first d moments of the distribution; with the linear kernel, KCP detects changes in the mean of the X_i .

From now on, we focus on the problem of detecting the changes of τ^* only, whether the kernel is characteristic or not.

2.4. Rewriting the empirical risk

It is convenient to see the images of the observations by the feature map as an element of \mathcal{H}^n . To this extent, we define $Y := (Y_1, \dots, Y_n)$, as well as $\mu^* := (\mu_1^*, \dots, \mu_n^*) \in \mathcal{H}^n$ and $\varepsilon := Y - \mu^* \in \mathcal{H}^n$. We identify the elements of \mathcal{H}^n with the set of applications $\{1, \dots, n\} \rightarrow \mathcal{H}$, naturally embedded with the inner product and norm given by

$$\forall x, y \in \mathcal{H}^n, \quad \langle x, y \rangle := \sum_{j=1}^n \langle x_j, y_j \rangle_{\mathcal{H}} \quad \text{and} \quad \|x\|^2 := \sum_{j=1}^n \|x_j\|_{\mathcal{H}}^2.$$

We now rewrite the empirical risk as a function of τ and Y . For any segmentation $\tau \in \mathcal{T}_n$, define F_τ the set of applications $\{1, \dots, n\} \rightarrow \mathcal{H}$ that are constant over the segments of τ . We see F_τ as a subspace of \mathcal{H}^n as a vector space. Take $f \in \mathcal{H}^n$, we define $\Pi_\tau f$ the orthogonal projection of f onto F_τ with respect to $\|\cdot\|$:

$$\Pi_\tau f \in \arg \min_{g \in F_\tau} \|f - g\|.$$

It is shown by Arlot, Celisse and Harchaoui [4] that for any $f \in \mathcal{H}^n$ and any $\ell \in \{1, \dots, D_\tau\}$,

$$\forall i \in \{\tau_{\ell-1} + 1, \dots, \tau_\ell\}, \quad (\Pi_\tau f)_i = \frac{1}{|\tau_\ell - \tau_{\ell-1}|} \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} f_j. \quad (2.5)$$

We are now able to write the empirical risk as

$$\widehat{\mathcal{R}}_n(\tau) = \frac{1}{n} \|Y - \widehat{\mu}_\tau\|^2, \quad (2.6)$$

where $\widehat{\mu}_\tau = \Pi_\tau Y$, following [27, 4].

2.5. Assumptions

A key ingredient of our analysis is the concentration of ε . Intuitively, the performance of KCP is better when ε concentrates strongly around its mean, since without noise we are just given the task to segment a piecewise-constant signal. It is thus natural to make assumptions on ε in order to obtain concentration results. We actually formulate assumptions on the kernel k , which translate automatically onto ε .

As done by Arlot, Celisse and Harchaoui [4], the main hypothesis used in our analysis is the following.

Assumption 1. A positive constant M exists such that

$$\forall i \in \{1, \dots, n\}, \quad k(X_i, X_i) \leq M^2 < +\infty \quad \text{a.s.}$$

If Assumption 1 holds true,

$$\forall i \in \{1, \dots, n\}, \quad \|Y_i\|_{\mathcal{H}} = \sqrt{k(X_i, X_i)} \leq M \quad \text{a.s.}$$

and Arlot, Celisse and Harchaoui [4] show that $\|\varepsilon_i\|_{\mathcal{H}} \leq 2M$ almost surely.

Assumption 1 is always satisfied for a large class of commonly used kernels, such as the Gaussian, Laplace and χ^2 kernels.

Note that Assumption 1 is weaker than assuming k to be bounded — that is, $k(x, x) \leq M$ for any $x \in \mathcal{X}$, which is equivalent to $k(x, x') \leq M$ for any $x, x' \in \mathcal{X}$ since k is positive definite. For instance, if $\mathcal{X} = \mathbb{R}^p$ and the data X_i are bounded almost surely, Assumption 1 holds true for the linear kernel and all polynomial kernels, which are not bounded on \mathbb{R}^p .

In the setting of Example 2.1, Assumption 1 holds true when

$$\forall j \in \{1, \dots, K\}, \quad k(x, x) \leq M^2 \quad \text{for } P_j\text{-a.e. } x \in \mathcal{X}.$$

It is sometimes possible to weaken Assumption 1 into a finite variance assumption.

Assumption 2. A positive constant $V < +\infty$ exists such that

$$\max_{1 \leq i \leq n} \mathbb{E} \left[\|\varepsilon_i\|_{\mathcal{H}}^2 \right] \leq V.$$

Since $v_i := \mathbb{E}[\|\varepsilon_i\|_{\mathcal{H}}^2] = \mathbb{E}[k(X_i, X_i)] - \|\mu_i^*\|_{\mathcal{H}}^2$, Assumption 2 holds true when

$$\forall i \in \{1, \dots, n\}, \quad \mathbb{E}[k(X_i, X_i)] \leq V.$$

As a consequence, Assumption 1 implies Assumption 2 with $V = M^2$. Note that Assumption 2 is satisfied for the polynomial kernel of order d provided that

$$\forall i \in \{1, \dots, n\}, \quad \mathbb{E} \left[\|X_i\|^{2d} \right] < +\infty.$$

In the setting of Example 2.1, Assumption 2 holds true with

$$V = \max_{1 \leq \ell \leq K+1} \mathbb{E}_{X \sim P_\ell} [k(X, X)],$$

provided this maximum is finite.

3. Theoretical guarantees for KCP

We are now able to state our main results. In Section 3.1, we state the main result of the paper, Theorem 3.1, which provides simple conditions under which KCP recovers the correct number of segments and localizes the true change-points with high probability, under the bounded kernel Assumption 1. Then, Section 3.2 details a few classical losses between segmentations which can be considered in addition to the one used in Theorem 3.1. Corollary 3.1 formulates a result on $\hat{\tau}$ in terms of the Frobenius loss. Finally, Section 3.3 states a partial result on KCP — requiring the number of change-points D_{τ^*} to be known — under the weaker Assumption 2.

3.1. Main result

We first need to define some quantities. The size of the smallest jump of μ^* in \mathcal{H} is

$$\underline{\Delta} := \min_{i / \mu_i^* \neq \mu_{i+1}^*} \|\mu_i^* - \mu_{i+1}^*\|_{\mathcal{H}}. \quad (3.1)$$

Intuitively, the higher $\underline{\Delta}$ is, the easier it is to detect the smallest jump with our procedure. The quantity $\|\mu_i^* - \mu_{i+1}^*\|_{\mathcal{H}}$ is often called the (population) maximum mean discrepancy [MMD, 24] between the distributions of X_i and X_{i+1} . In the scalar setting (with the linear kernel), the ratio $\underline{\Delta}/\sigma$ (where σ^2 is the variance of the noise) is called the *signal-to-noise ratio* [7] and is often used as a measure of the magnitude of a change in the signal. In Example 2.1,

$$\underline{\Delta} = \min_{1 \leq j \leq K} \|\mu_{P_j}^* - \mu_{P_{j+1}}^*\|_{\mathcal{H}}$$

where $\mu_{P_j}^*$ denotes the (Bochner) expectation of $\Phi(X)$ when $X \sim P_j$.

For any $\tau \in \mathcal{T}_n$, we denote the (normalized) sizes of its smallest and of its largest segment by

$$\underline{\Lambda}_\tau := \frac{1}{n} \min_{1 \leq \ell \leq D_\tau} |\tau_\ell - \tau_{\ell-1}| \quad \text{and} \quad \bar{\Lambda}_\tau := \frac{1}{n} \max_{1 \leq \ell \leq D_\tau} |\tau_\ell - \tau_{\ell-1}|. \quad (3.2)$$

It should be clear that the smaller $\underline{\Lambda}_{\tau^*}$ is, the harder it is to detect the segment that achieves the minimum in equation (3.2). For instance, in the particular case of Example 2.1,

$$\underline{\Lambda}_{\tau^*} \xrightarrow{n \rightarrow +\infty} \min_{0 \leq j \leq K} |b_{j+1} - b_j| \quad \text{and} \quad \bar{\Lambda}_{\tau^*} \xrightarrow{n \rightarrow +\infty} \max_{0 \leq j \leq K} |b_{j+1} - b_j|.$$

For any τ^1 and $\tau^2 \in \mathcal{T}_n$, we define

$$d_\infty^{(1)}(\tau^1, \tau^2) := \max_{1 \leq i \leq D_{\tau^1}-1} \left\{ \min_{1 \leq j \leq D_{\tau^2}-1} |\tau_i^1 - \tau_j^2| \right\},$$

which is a loss function (a measure of dissimilarity) between the segmentations τ^1 and τ^2 . Note that $d_\infty^{(1)}$ is not a distance; other possible losses between segmentations and their relationship with $d_\infty^{(1)}$ are discussed in Section 3.2.

Theorem 3.1. *Suppose that Assumption 1 holds true. For any $y > 0$, an event Ω of probability at least $1 - e^{-y}$ exists on which the following holds true. For any $C > 0$, let $\hat{\tau}$ be defined as in Eq. (2.2) with pen defined by Eq. (2.3). Set*

$$C_{\min} := \frac{74}{3}(D_{\tau^*} + 1)(y + \log n + 1) \quad \text{and} \quad C_{\max} := \frac{\underline{\Delta}^2}{M^2} \frac{\underline{\Lambda}_{\tau^*}}{6D_{\tau^*}} n.$$

Then, if

$$C_{\min} < C < C_{\max}, \quad (3.3)$$

on Ω , we have

$$D_{\hat{\tau}} = D_{\tau^*} \quad \text{and} \quad \frac{1}{n} d_\infty^{(1)}(\tau^*, \hat{\tau}) \leq v_1(y) := \frac{148 D_{\tau^*} M^2}{\underline{\Delta}^2} \cdot \frac{y + \log n + 1}{n}.$$

We delay the proof of Theorem 3.1 to Section 6.4. Some remarks follow.

Theorem 3.1 is a non-asymptotic result: it is valid for any $n \geq 1$ and there is nothing hidden in $o(1)$ remainder terms. The latter point is crucial for complex data — for instance, $\mathcal{X} = \mathbb{R}^p$ with $p > n$ — since in this case, assuming \mathcal{X} fixed while $n \rightarrow +\infty$ is not realistic.

Nevertheless, it is useful to write down what Theorem 3.1 becomes in the asymptotic setting of Example 2.1. As previously noticed, D_{τ^*} , $\underline{\Delta}_{\tau^*}$, $\underline{\Delta}^2$ and M^2 then converge to positive constants as $n \rightarrow +\infty$. Therefore, C_{\min} is of order $\log(n)$, C_{\max} is of order n and we always have $C_{\min} < C_{\max}$ for n large enough. The upper bound on C matches classical asymptotic conditions for variable selection [47]. The necessity of taking C of order at least $\log(n)$ is shown by Birgé and Massart [11] in a variable selection setting, which includes change-point detection as a particular example; Birgé and Massart [11], Abramovich et al. [2] provide several arguments for the optimality of taking a constant C of order $\log(n)$. When C satisfies (3.3), the result of Theorem 3.1 implies that $\mathbb{P}(D_{\hat{\tau}} = D_{\tau^*}) \rightarrow 1$. For the linear kernel in \mathbb{R}^d , this is a well-known result when the distribution of the X_i changes only through its mean. The first result dates back to Yao [54, Section 2] for a Gaussian noise, later extended by Liu, Wu and Zidek [42] and Bai and Perron [6, Section 3.1] under mixingale hypothesis on the error, and Lavielle and Moulines [37] under very mild assumptions satisfied for a large family of zero-mean processes [for the precise statement of the hypothesis, see 37, Section 2.1]. Theorem 3.1 also shows that the normalized estimated change-points of $\hat{\tau}$ converge towards the normalized true change-points at speed at least $\log(n)/n$.

Up to a logarithmic factor, this speed matches the minimax lower bound n^{-1} which has been obtained previously for various change-point procedures [32, 12, 33, for instance] including least-squares [37], assuming that $\underline{\Delta}_{\tau^*} \geq \kappa > 0$. When $D_{\tau^*} \geq 3$ and the assumption on $\underline{\Delta}_{\tau^*}$ is removed —that is, segments of length much smaller than n are allowed, which is compatible with Theorem 3.1 since it is non-asymptotic—, Brunel [14, Theorem 6] shows a minimax lower bound of order $\log(n)/n$. Therefore, in this setting, KCP achieves the minimax rate. We do not know whether KCP remains minimax optimal (without the log factor) under the assumption $\underline{\Delta}_{\tau^*} \geq \kappa > 0$.

Note finally that KCP also performs well for finite samples, according to the simulation experiments of Arlot, Celisse and Harchaoui [4].

Theorem 3.1 emphasizes the key role of $\underline{\Delta}^2/M^2$, which can be seen as a generalization of the signal-to-noise ratio, for the change-point detection performance of KCP. The larger is this ratio, the easier it is to have Eq. (3.3) satisfied and the smaller is $v_1(y)$. This suggests to choose k (theoretically at least) by maximizing $\underline{\Delta}^2/M^2$, as we discuss in Section 5. Note that $\underline{\Delta}^2/M^2$ is invariant by a rescaling of k , hence the result of Theorem 3.1 is unchanged when k is rescaled.

The hypothesis in Eq. (3.3) is actually three-fold. First, we use that $C > C_{\min}$ to get $D_{\hat{\tau}} \leq D_{\tau^*}$. We have to assume C large enough since a too small penalty leads to selecting (with KCP or any other penalized least-squares procedure) the segmentation with n segments, that is $D_{\hat{\tau}} = n$. Second, $C < C_{\max}$ is used to

get $D_{\hat{\tau}} \geq D_{\tau^*}$. Such an assumption is required since taking a penalty function too large in Eq. (2.2) would result in selecting the segmentation with only one segment, that is, $D_{\hat{\tau}} = 1$. Third, C_{\max} has to be greater than C_{\min} for providing a non-empty interval of possible values for C . This inequality is also used in the proof of the upper bound on $d_{\infty}^{(1)}(\tau^*, \hat{\tau})$ when we already know that $D_{\hat{\tau}} = D_{\tau^*}$. In Example 2.1, the $C_{\min} < C_{\max}$ hypothesis translates into $\underline{\Delta}_{\tau^*} \succ \log(n)/n$. That is, the size of the smallest segment has to be of order $\log n/n$. This is known to be a necessary condition to obtain the minimax rate in multiple change-point detection [14, section 2].

Theorem 3.1 helps choosing C , which is a key parameter of KCP, as in any penalized model selection procedure. However, in practice, we do not recommend to directly use equation (3.3) for choosing C for two reasons: C_{\min}, C_{\max} depend on unknown quantities $D_{\tau^*}, \underline{\Delta}_{\tau^*}, \underline{\Delta}$, and the exact values of the constants in C_{\min}, C_{\max} might be pessimistic compared to what we can observe from simulation experiments. We rather suggest to use a data-driven method for choosing C , see Section 5.

If we know D_{τ^*} , we can replace $\hat{\tau}$ by

$$\hat{\tau}(D_{\tau^*}) \in \arg \min_{\tau \in \mathcal{T}_n^{D_{\tau^*}}} \{\hat{\mathcal{R}}_n(\tau)\}.$$

Then, assuming that $\underline{\Delta}_{\tau^*} > v_1(y)$ — which is weaker than assuming $C_{\min} < C_{\max}$ —, the proof of Theorem 3.1 shows that, on Ω , we have

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}(D_{\tau^*})) \leq v_1(y).$$

3.2. Loss functions between segmentations

Theorem 3.1 shows that $\hat{\tau}$ is close to τ^* in terms of $d_{\infty}^{(1)}$. Several other loss functions (measures of dissimilarity) can be defined between segmentations [29]. We here consider a few of them, which are often used or natural for the change-point problem.

Let us first consider losses related to the Hausdorff distance. For any τ^1 and $\tau^2 \in \mathcal{T}_n$, we define

$$\begin{aligned} d_{\infty}^{(1)}(\tau^1, \tau^2) &:= \max_{1 \leq i \leq D_{\tau^1}-1} \left\{ \min_{1 \leq j \leq D_{\tau^2}-1} |\tau_i^1 - \tau_j^2| \right\} \\ d_{\infty}^{(2)}(\tau^1, \tau^2) &:= \max_{1 \leq i \leq D_{\tau^1}-1} \left\{ \min_{0 \leq j \leq D_{\tau^2}} |\tau_i^1 - \tau_j^2| \right\} \\ d_H^{(i)}(\tau^1, \tau^2) &:= \max \{ d_{\infty}^{(i)}(\tau^1, \tau^2), d_{\infty}^{(i)}(\tau^2, \tau^1) \} \quad \text{for } i \in \{1, 2\}. \end{aligned}$$

Whenever $D_{\tau^1} = D_{\tau^2}$, we define

$$d_{\infty}^{(3)}(\tau^1, \tau^2) := \max_{1 \leq i \leq D_{\tau^1}-1} |\tau_i^1 - \tau_i^2|.$$

Note that $d_\infty^{(3)}$ is symmetric thus there is no need to define $d_H^{(3)}$. One could also define $d_H^{(1)}$ as the *Hausdorff distance* between the subsets $\{\tau_1^1, \dots, \tau_{D_{\tau^1}-1}^1\}$ and $\{\tau_1^2, \dots, \tau_{D_{\tau^2}-1}^2\}$ with respect to the distance $\delta(x, y) = |x - y|$ on \mathbb{R} . These definitions are illustrated by Figure 3.

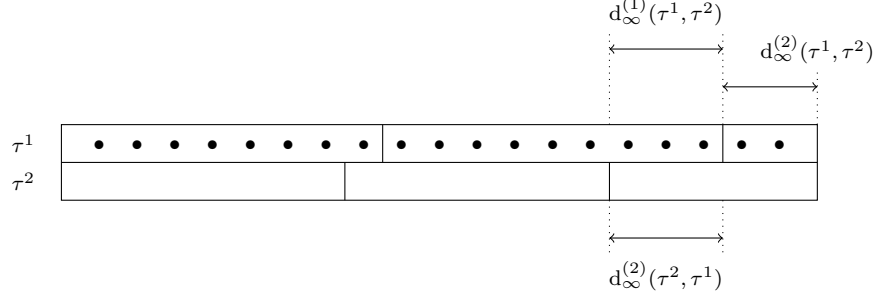


FIG 3. Illustration of the definition of $d_\infty^{(i)}$, with $n = 19$, $\tau^1 = [0, 8, 17, 19]$ and $\tau^2 = [0, 7, 14, 19]$. In this example, $D_{\tau^1} = D_{\tau^2} = 3$. We can compute $d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^2, \tau^1) = d_\infty^{(2)}(\tau^2, \tau^1) = d_\infty^{(2)}(\tau^1, \tau^2) = 3$ and $d_\infty^{(3)}(\tau^1, \tau^2) = 2$.

Interestingly, all these loss functions coincide whenever $n^{-1}d_\infty^{(1)}(\tau^1, \tau^2)$ is small enough. The following lemma makes this claim rigorous.

Lemma 3.1. *We have the following two properties.*

(i) *For any $\tau^1, \tau^2 \in \mathcal{T}_n$ such that*

$$\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \frac{1}{2} \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\},$$

we have $D_{\tau^1} = D_{\tau^2}$ and

$$d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(2)}(\tau^1, \tau^2) = d_\infty^{(3)}(\tau^1, \tau^2) = d_H^{(1)}(\tau^1, \tau^2) = d_H^{(2)}(\tau^1, \tau^2).$$

(ii) *For any $\tau^1, \tau^2 \in \mathcal{T}_n$ such that*

$$D_{\tau^1} = D_{\tau^2} \quad \text{and} \quad \frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \frac{\underline{\Lambda}_{\tau^1}}{2},$$

we have

$$d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^2, \tau^1) = d_H^{(1)}(\tau^1, \tau^2).$$

Lemma 3.1 is proved in Section B.1. As a direct application of Lemma 3.1 we see that the statement of Theorem 3.1 holds true with $d_\infty^{(1)}$ replaced by *any* of the loss functions that we defined above, at least for n large enough.

Another loss between segmentations is the *Frobenius* loss [35], which is defined as follows. For any $\tau^1, \tau^2 \in \mathcal{T}_n$,

$$d_F(\tau^1, \tau^2) := \|\Pi_{\tau^1} - \Pi_{\tau^2}\|_F,$$

where Π_τ is the orthogonal projection onto F_τ , as defined in Section 2.4, and $\|\cdot\|_F$ denotes the Frobenius norm of a matrix:

$$\forall A \in \mathbb{R}^{N \times M}, \quad \|A\|_F^2 := \sum_{i=1}^N \sum_{j=1}^M A_{ij}^2.$$

A closed-form formula for d_F can be derived from the matrix representation of Π_τ that is given by (2.5): for any $i, j \in \{1, \dots, n\}$,

$$(\Pi_\tau)_{i,j} = \begin{cases} \frac{1}{|\lambda|} & \text{if } i \text{ and } j \text{ belong to the same segment } \lambda \text{ of } \tau \\ 0 & \text{otherwise.} \end{cases}$$

An interesting feature of the Frobenius loss is that it is smaller than one only when τ^1 and τ^2 have the same number of segments, whereas Hausdorff distances can be small with very different numbers of segments. Indeed, we prove in Section B.2 that

$$|D_{\tau^1} - D_{\tau^2}| \leq d_F(\tau^1, \tau^2)^2 \leq D_{\tau^1} + D_{\tau^2}. \quad (3.4)$$

The next proposition shows that there is an equivalence (up to constants) between the Hausdorff and Frobenius losses between segmentations, provided that they are close enough.

Proposition 3.1. *Suppose that $D_{\tau^1} = D_{\tau^2}$ and $\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \underline{\Lambda}_{\tau^1}/2$, then*

$$(d_F(\tau^1, \tau^2))^2 \leq \frac{12D_{\tau^1}}{\underline{\Lambda}_{\tau^1}} \frac{1}{n} d_\infty^{(1)}(\tau^1, \tau^2).$$

If in addition $\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \underline{\Lambda}_{\tau^1}/3$, then

$$\frac{2}{3\underline{\Lambda}_{\tau^1}} \frac{1}{n} d_\infty^{(1)}(\tau^1, \tau^2) \leq (d_F(\tau^1, \tau^2))^2.$$

Proposition 3.1 was first stated and proved by [35, Theorem B.2]. We prove it in Section B.2 for completeness.

As a corollary of Theorem 3.1 and Proposition 3.1, we get the following guarantee on the Frobenius loss between τ^\star and the segmentation $\hat{\tau}$ estimated by KCP.

Corollary 3.1. *Under the assumptions of Theorem 3.1, on the event Ω defined by Theorem 3.1, for any $\hat{\tau}$ satisfying (2.2) with pen defined by (2.3), we have:*

$$d_F(\tau^\star, \hat{\tau}) \leq \frac{43D_{\tau^\star}}{\sqrt{\underline{\Lambda}_{\tau^\star}}} \cdot \frac{M}{\underline{\Delta}} \sqrt{\frac{y + \log n + 1}{n}}.$$

Note that Corollary 3.1 gives a better result (at least for large n) than the obvious bound

$$d_F(\tau^\star, \hat{\tau}) \leq D_{\tau^\star} + D_{\hat{\tau}} - 2.$$

Proof. On the event Ω , we have $\frac{1}{n}d_{\infty}^{(1)}(\tau^*, \hat{\tau}) < \underline{\Delta}_{\tau^*}/(D_{\tau^*} + 1)$ and $D_{\tau^*} = D_{\hat{\tau}}$. Therefore, according to Proposition 3.1,

$$(d_F(\tau^*, \hat{\tau}))^2 \leq \frac{12D_{\tau^*}}{\underline{\Delta}_{\tau^*}} \frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}) \leq \frac{1776D_{\tau^*}^2(y + \log n + 1)}{n\underline{\Delta}_{\tau^*}} \cdot \frac{M^2}{\underline{\Delta}^2}.$$

□

Up to this point, we assessed the quality of the segmentation τ by considering the proximity of τ with τ^* . Another natural idea is to measure the distance between μ^* and μ_{τ}^* in \mathcal{H}^n . It is closely related to the oracle inequality proved by Arlot, Celisse and Harchaoui [4], which implies an upper bound on $\|\mu^* - \hat{\mu}_{\hat{\tau}}\|^2$. We can also observe that there is a simple relationship between $\|\mu^* - \mu_{\tau}^*\|^2$ and the Frobenius distance between τ and τ^* . Indeed,

$$\|\mu^* - \mu_{\tau}^*\|^2 = \|(\Pi_{\tau^*} - \Pi_{\tau})\mu^*\|^2 \leq \|\Pi_{\tau^*} - \Pi_{\tau}\|_2^2 \|\mu^*\|^2 \leq (d_F(\tau^*, \hat{\tau}))^2 \|\mu^*\|^2. \quad (3.5)$$

Equation (6.9) in the proof of Theorem 3.1 shows that on Ω , under the assumptions of Theorem 3.1,

$$\|\mu^* - \mu_{\hat{\tau}}^*\|^2 \leq 74(y + \log(n) + 1)D_{\tau^*}M^2$$

which is slightly better (but similar) to what Corollary 3.1, equation (3.5) and the bound $\|\mu^*\|^2 \leq M^2n$ imply.

3.3. Extension to the finite variance case

Theorem 3.1 is valid under a boundedness assumption (Assumption 1). What happens under the weaker Assumption 2? As a first step, we provide a result for

$$\hat{\tau}(D_{\tau^*}, \delta_n) \in \arg \min_{\tau \in \mathcal{T}_n^{D_{\tau^*}} / \underline{\Delta}_{\tau} \geq \delta_n} \{\hat{\mathcal{R}}_n(\tau)\} \quad (3.6)$$

for some $\delta_n > 0$. In other words, we restrict our search to segmentations τ of the correct size — hence D_{τ^*} must be known *a priori* — and having no segment with less than $n\delta_n$ observations. We discuss how to relax this restriction right after the statement of Theorem 3.2. Note that the dynamic programming algorithm of Harchaoui and Cappé [27] can be used for computing $\hat{\tau}(D_{\tau^*}, \delta_n)$ efficiently.

Similarly to $\underline{\Delta}$, we define $\overline{\Delta} := \max_i \|\mu_i^* - \mu_{i+1}^*\|_{\mathcal{H}}$.

Theorem 3.2. *Suppose that Assumption 2 holds true. For any $\delta_n, y > 0$, define:*

$$v_2(y, \delta_n) := 24(D_{\tau^*})^2 \frac{\overline{\Delta}\sqrt{V}}{\underline{\Delta}^2} \frac{y}{\sqrt{n}} + 8D_{\tau^*} \frac{V}{\underline{\Delta}^2} \frac{y^2}{n\delta_n}.$$

For any $y > 0$, an event Ω_2 exists such that

$$\mathbb{P}(\Omega_2) \geq 1 - \frac{1}{y^2}$$

and, on Ω_2 , we have the following: for any $\delta_n \in (0, \underline{\Delta}_{\tau^*}]$ and any $\hat{\tau}(D_{\tau^*}, \delta_n)$ satisfying Eq. (3.6), if $v_2(y, \delta_n) \leq \underline{\Delta}_{\tau^*}$,

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}(D_{\tau^*}, \delta_n)) \leq v_2(y, \delta_n). \quad (3.7)$$

We postpone the proof of Theorem 3.2 to Section 6.5. Let us make a few remarks.

As for Theorem 3.1, our result is non-asymptotic. However, it is interesting to write it down in the setting of Example 2.1. If n goes to infinity, then the assumption $\underline{\Delta}_{\tau^*} \geq \delta_n$ is satisfied whenever $\delta_n \rightarrow 0$. If we furthermore require that $n\delta_n \rightarrow \infty$, then Eq. (3.7) implies that

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}(D_{\tau^*}, \delta_n)) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} 0,$$

by taking a well-chosen y of order $\sqrt{n} + \sqrt{n\delta_n}$. In the particular case of the linear kernel, this result is known under various hypothesis [37, for instance]; it is new for a general kernel.

More precisely, if we take $\delta_n = n^{-1/2}$, Theorem 3.2 implies that

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}(D_{\tau^*}, n^{-1/2}))$$

goes to zero at least as fast as ℓ_n/\sqrt{n} , where $(\ell_n)_{n \geq 1}$ is any sequence tending to infinity, for instance $\ell_n = \log(n)$. This speed seems suboptimal compared to previous results [37, for instance] — which do not consider the case of a general kernel —, but we have not been able to prove tight enough deviation bounds for getting the localization rate $\log(n)/n$ under Assumption 2.

How does Theorem 3.2 compares to Theorem 3.1? First, as noticed by Remark 6.4 in Section 6.4, the result of Theorem 3.1 also holds true for $\hat{\tau}(D_{\tau^*}, \delta_n)$ as long as $\underline{\Delta}_{\tau^*} \geq \delta_n$. Second, $v_1(y)$ is usually smaller than $v_2(y, \delta_n)$ — its order of magnitude is smaller when $n \rightarrow +\infty$ —, and the lower bound on the probability of Ω is better than the one for Ω_2 . There is no surprise here: the stronger Assumption 1 helps us proving a stronger result for $\hat{\tau}(D_{\tau^*}, \delta_n)$. Nevertheless, these only are upper bounds, so we do not know whether the performance of $\hat{\tau}(D_{\tau^*}, \delta_n)$ actually changes much depending on the noise assumption. For instance, as already noticed, we do not believe that the localization speed $\log(n)/n$ requires a boundedness assumption; in particular cases at least, it has been obtained for unbounded data [37, 12].

The dependency in k of the speed of convergence of $\hat{\tau}(D_{\tau^*}, \delta_n)$ is slightly less clear than in Theorem 3.1. The signal-to-noise ratio appears through $\underline{\Delta}^2/V$, as expected, but the size $\underline{\Delta}$ of the largest true jump also appears in v_2 . At the very least, it is clear that $\underline{\Delta}^2/V$ should not be too small.

As noted by Lavielle and Moulines [37], it may be possible to get rid of the minimal segment length δ_n , either by imposing stronger conditions on ε — which are not met in our setting — or by constraining the values of $\hat{\mu}$ to lie in a compact subset $\Theta \subset \mathcal{H}^{D_{\tau^*}+1}$.

4. Numerical simulations

One consequence of our main result, Theorem 3.1, is that for a bounded kernel, the KCP procedure is consistent in the asymptotic setting presented in Example 2.1. We now illustrate this fact by a simulation study.

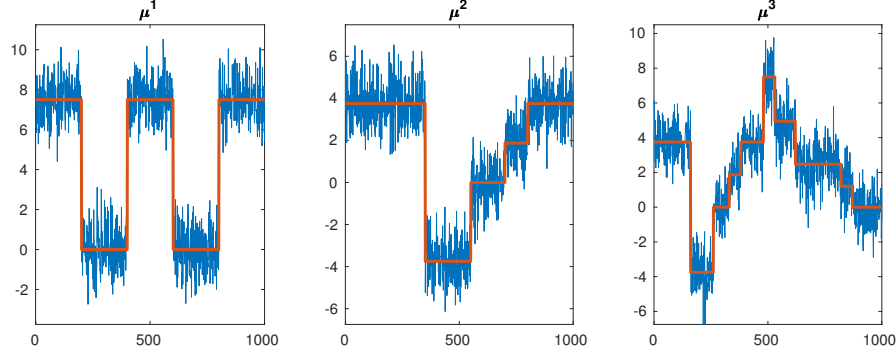


FIG 4. In red, the three piecewise constant functions used in the simulations of Section 4. In blue, a noisy version of these functions. Both μ^1 and μ^2 have 4 jumps; μ^3 has 9 jumps.

Detecting changes in the mean with the Gaussian kernel Let us consider the archetypic change-point detection problem —finding changes in the mean of a sequence of independent random variables— and show how these changes are localized more precisely when more data are available.

We define three functions $\mu^m : [0, 1] \rightarrow \mathbb{R}$, $1 \leq m \leq 3$, previously used by Arlot and Celisse [3], which cover a variety of situations (see Fig. 4). For each $m \in \{1, 2, 3\}$ and several values of n between 10^2 and 10^3 , we repeat 10^3 times the following:

- Sample n independent Gaussian random variables $g_i \sim \mathcal{N}(0, 1)$;
- Set $X_i = \mu^m(i/n) + g_i$ —Fig. 4 shows one sample for each $m \in \{1, 2, 3\}$;
- Perform KCP with Gaussian kernel and linear penalty on X_1, \dots, X_n ; the penalty constant is chosen as indicated in Section 5, the bandwidth is set to 0.1, and the maximum number of change-points is set to 30;
- Compute $d_H^{(2)}(\tau^*, \hat{\tau}_n)$.

The results are collected in Fig. 5, where each graph corresponds to a regression function μ^m . We represent in logarithmic scale the mean distance between the true segmentation and the estimated segmentation for each value of n . The error bars are $\pm \hat{\sigma}/\sqrt{N}$, where $\hat{\sigma}$ is the empirical standard deviation over $N = 10^3$ repetitions. We want to emphasize that, though these experiments illustrate our main result Theorem 3.1, they are carried out in a slightly different setting since the penalty constant C is not chosen according to equation (3.3), but using the dimension jump heuristic [8].

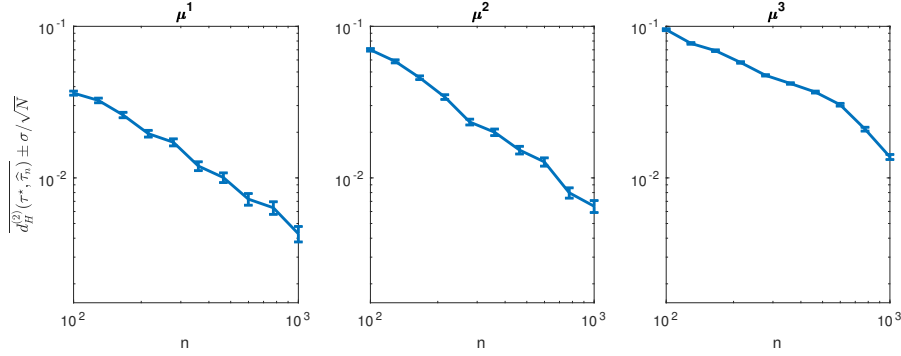


FIG 5. Convergence of $\frac{1}{n}d_H^{(2)}(\tau^*, \hat{\tau}_n)$ towards 0 when the number of data points n is increasing. A linear regression between $\log n$ and $\frac{1}{n}d_H^{(2)}(\tau^*, \hat{\tau}_n)$ for $n \geq 300$ yields slope estimates -0.97 , -1.04 and -1.00 , respectively.

The three segmentation problems considered here are quite different in nature, but all lead to a linear convergence rate (slopes close to -1 on the graphs of Figure 5) with different constants (different values for the intercept on the graphs of Figure 5). Recall that Theorem 3.1 combined with Lemma 3.1 states that, with high probability,

$$\frac{1}{n}d_H^{(2)}(\tau^*, \hat{\tau}_n) \lesssim \tilde{v}_1 = \frac{D_{\tau^*} M^2}{\underline{\Delta}^2} \cdot \frac{\log n}{n}.$$

Hence, whenever D_{τ^*} , $\underline{\Delta}$ and M are fixed, $\frac{1}{n}d_H^{(2)}(\tau^*, \hat{\tau}_n)$ converges to 0 at rate at least $\log n/n$ when the number of data points increases. In our experimental setting, these quantities are fixed, and the observed convergence rate matches our theoretical upper bound. The performance of KCP still depends on the regression function μ^m experimentally, by a constant multiplicative factor, like the theoretical bound \tilde{v}_1 .

Detecting changes in the number of modes Let us now consider data $X_1, \dots, X_n \in \mathbb{R}$ whose distribution vary only through the number of modes. Can we accurately detect such changes with the KCP procedure? The data are generated according to the following process for several n :

- Set $\tau_1^* = \lfloor n/3 \rfloor$ and $\tau_2^* = \lfloor 2n/3 \rfloor$;
- Draw $X_1, \dots, X_{\tau_1^*}, X_{\tau_2^*+1}, \dots, X_n$ according to a standard Gaussian distribution, and $X_{\tau_1^*+1}, \dots, X_{\tau_2^*}$ according to a $(1/2, 1/2)$ -mixture of Gaussian distributions $\mathcal{N}(\delta, 1 - \delta^2)$ and $\mathcal{N}(-\delta, 1 - \delta^2)$, with $\delta = 0.999$; the X_i are independent.

We test KCP with various kernels assuming that the number of change-points ($D_{\tau^*} = 3$) is known; this simplification avoids possible artifacts linked to the choice of the penalty constant. Results are shown on Figure 6. The X_i all have

zero mean and unit variance, hence a classical penalized least-squares procedure —KCP with the linear kernel— is expected to detect poorly the changes in the distribution of the X_i , as confirmed by Figure 6 (for instance, according to the right panel, it is not consistent). On the contrary, a Gaussian kernel with well-chosen bandwidth yields much better performance according to the middle and right panels of Figure 6 (with a rate of order $1/n$).

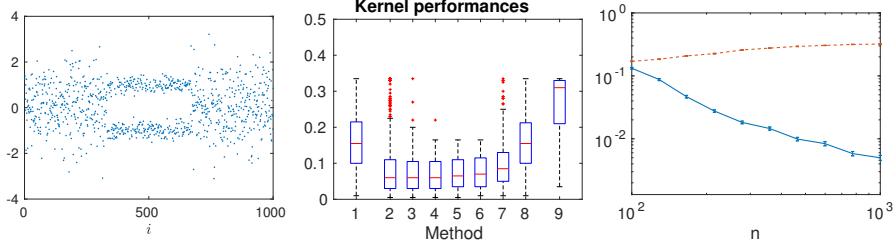


FIG 6. **Left:** One sample X_1, \dots, X_n for $n = 10^3$. **Middle:** Performance of KCP with various kernels ($n = 200$). Methods 1 to 8: Gaussian kernel with bandwidth set via the median heuristic (method 1), or fixed equal to 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1 (methods 2, ..., 9, respectively). Method 9: linear kernel. **Right:** Estimated values of $n^{-1}d_H^{(2)}(\tau^*, \hat{\tau}_n)$ vs. n in log scale, for KCP with a Gaussian kernel with bandwidth 0.01 (blue solid line; estimated slope -1.05) and with the linear kernel (red dashed line; estimated slope 0.16).

5. Discussion

Before proving our main results, let us discuss some of their consequences regarding the KCP procedure.

Fully non-parametric consistent change-point detection We have proved that for any kernel satisfying some reasonably mild hypotheses, the KCP procedure outputs a segmentation closeby the true segmentation with high probability.

An important particular example is the “asymptotic setting” of Example 2.1, where we have a fixed true segmentation τ^* and fixed distributions P_1, \dots, P_{K+1} from which more and more points are sampled. How fast can KCP recover τ^* , without any prior information on the number of segments D_{τ^*} or on the distributions P_1, \dots, P_{K+1} ?

Let us take a bounded characteristic kernel — for instance the Gaussian or the Laplace kernel if $\mathcal{X} = \mathbb{R}^d$ —, so that Assumption 1 holds true. Then, Theorem 3.1 shows that KCP detects consistently all changes in the distribution of the X_i , and localizes them at speed $\log(n)/n$. This speed also depends on the adequation between the kernel k and the differences between the P_j , through the ratio $\underline{\Delta}^2/M^2$. Obtaining such a fully non-parametric result for multiple change-points with a general set \mathcal{X} — we only need to know a bounded characteristic kernel on \mathcal{X} — has never been obtained before. To the best of our knowledge, non-parametric consistency results for the detection of arbitrary changes in the

distribution of the data have only been obtained for real-valued data [56] or for the case of a single change-point [15, 13].

Choice of k An important question remains: how to choose the kernel k ? In Theorem 3.1, k only appears through the “signal-to-noise ratio” $\underline{\Delta}^2/M^2$, leading to better theoretical guarantees when this signal-to-noise ratio is larger: a larger value for C_{\max} and a smaller bound v_1 on $d_{\infty}^{(1)}(\tau^*, \hat{\tau})$. Therefore, a simple strategy for choosing the kernel is to pick k that maximizes $\underline{\Delta}^2/M^2$, at least among a family of kernels, for instance Gaussian kernels. This first idea requires to know the distributions of the X_i , or at least to have prior information on them. Interestingly, when the change-points locations are known, $\underline{\Delta}^2$ corresponds to the maximum mean discrepancy [MMD, 24] between the distributions of the X_i over contiguous segments. In this particular setting, it is feasible to estimate and to maximize $\underline{\Delta}^2$ with respect to the kernel k , as done by Gretton et al. [25]. An interesting future development would be to build an estimator of $\underline{\Delta}^2$ without knowing the change-point locations and to maximize this estimator with respect to the kernel k . We refer to Arlot, Celisse and Harchaoui [4, section 7.2] for a complementary discussion about the choice of k for KCP.

Choice of C Another important parameter of the KCP procedure is the constant C that appears in the penalty function. As mentioned below Theorem 3.1, our theoretical guarantees provide some guidelines for choosing C , but these are not sufficient to choose precisely C in practice. We recommend to follow the advice of [4, section 6.2] on this point, which is to choose C from data with the “slope heuristic” [8].

Modularity of the proofs and possible extensions Finally, we would like to emphasize what we believe to be an important contribution of this paper. The structure of the proofs of Theorems 3.1 and 3.2 — which follow the same strategy — is modular, so that one can easily adapt it to different sets of assumptions.

Our proof strategy is not fully new, since it is similar to the one of almost all previous papers analyzing the consistency of least-squares change-point detection procedures. In particular, we adapted some ideas of the proofs of Lavielle and Moulines [37] to the Hilbert space setting. Nevertheless, these papers formulate their main results in asymptotic terms, which can be seen as a limitation — especially when n is small or \mathcal{X} is of large dimension. Another approach is the one of Lebarbier [39], Comte and Rozenholc [17], Arlot, Celisse and Harchaoui [4] where non-asymptotic oracle inequalities — using concentration inequalities and following the model selection results of Birgé and Massart [10] — are provided as theoretical guarantees on some penalized least-squares change-point procedures. Up to now, these two approaches seemed difficult to combine. The proofs of Theorems 3.1 and 3.2 show how they can be reconciled, which allows us to mix their strengths.

Indeed, the assumptions on the distributions of the X_i —Assumptions 1 and 2— are only used for proving bounds on two quantities —a linear term L_τ and a quadratic term Q_τ —, uniformly over $\tau \in \mathcal{T}_n$. Under Assumption 1, this is done thanks to concentration inequalities (Lemmas 6.7 and 6.8) which have been proved first by Arlot, Celisse and Harchaoui [4] in order to get an oracle inequality. Under Assumption 2, this is done by generalizing the method of Lavielle and Moulines [37] to Hilbert-space valued data, through two deterministic bounds (Lemmas 6.5 and 6.6) and a deviation inequality for

$$M_n := \max_{1 \leq k \leq n} \left\| \sum_{j=1}^k \varepsilon_j \right\|_{\mathcal{H}}$$

(Lemma 6.10). The rest of the proofs does not use anything about the distribution of X_1, \dots, X_n .

As a consequence, if one can generalize these bounds to another setting, a straightforward consequence is that a result similar to Theorem 3.1 or 3.2 holds true for the KCP procedure in this new setting. In particular, this could be used for dealing with the case of dependent data X_1, \dots, X_n . We could also consider an intermediate assumption between Assumption 2 and Assumption 1, of the form:

$$\max_{1 \leq i \leq n} \mathbb{E}[k(X_i, X_i)^\alpha] \leq B_\alpha < +\infty$$

for some $\alpha \in (1, +\infty)$.

6. Proofs

Let us start by describing our general strategy for proving our main results. Our goal is to build a large probability event on which any $\hat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n} \text{crit}(\tau)$ belongs to some subset \mathcal{E} of \mathcal{T}_n . For proving this, we use the key fact that $\text{crit}(\tau^*) \geq \text{crit}(\hat{\tau})$, together with a lower bound on $\text{crit}(\tau)$ holding simultaneously for all $\tau \in \mathcal{T}_n$ —hence for $\tau = \hat{\tau}$.

In order to get such a lower bound on the empirical penalized criterion, we start by decomposing it in Section 6.1 into terms that are simpler to control individually: two random terms — a linear function of ε and a quadratic function of ε —, and two deterministic terms — the approximation error and the penalty. Then, we control these terms thanks to deterministic bounds (Section 6.2) and deviation/concentration inequalities (Section 6.3). Finally, we prove Theorem 3.1 in Section 6.4 and Theorem 3.2 in Section 6.5.

6.1. Decomposition of the empirical risk

The first step in the proofs of Theorems 3.1 and 3.2 is to decompose the empirical risk (2.6).

Lemma 6.1. *Let $\tau \in \mathcal{T}_n$ be a segmentation. Define $\mu_\tau^\star = \Pi_\tau \mu^\star$. Then we can write*

$$n\widehat{\mathcal{R}}_n(\tau) = \|Y - \widehat{\mu}_\tau\|^2 = \|\mu^\star - \mu_\tau^\star\|^2 + 2\langle \mu^\star - \mu_\tau^\star, \varepsilon \rangle - \|\Pi_\tau \varepsilon\|^2 + \|\varepsilon\|^2. \quad (6.1)$$

Proof. First, recall that $\widehat{\mu}_\tau = \Pi_\tau Y$ and that $Y = \mu^\star + \varepsilon$, hence

$$\begin{aligned} \|Y - \widehat{\mu}_\tau\|^2 &= \|Y - \Pi_\tau Y\|^2 \\ &= \|\mu^\star + \varepsilon - \Pi_\tau(\mu^\star + \varepsilon)\|^2 \\ &= \|\mu^\star - \Pi_\tau \mu^\star\|^2 + \|\varepsilon - \Pi_\tau \varepsilon\|^2 + 2\langle \mu^\star - \Pi_\tau \mu^\star, \varepsilon - \Pi_\tau \varepsilon \rangle. \end{aligned}$$

Since Π_τ is an orthogonal projection,

$$\begin{aligned} \|Y - \widehat{\mu}_\tau\|^2 &= \|\mu^\star - \mu_\tau^\star\|^2 + \|\varepsilon\|^2 - 2\langle \varepsilon, \Pi_\tau \varepsilon \rangle + \|\Pi_\tau \varepsilon\|^2 + 2\langle (\text{Id} - \Pi_\tau)\mu^\star, \varepsilon \rangle \\ &= \|\mu^\star - \mu_\tau^\star\|^2 + \|\varepsilon\|^2 - \|\Pi_\tau \varepsilon\|^2 + 2\langle (\text{Id} - \Pi_\tau)\mu^\star, \varepsilon \rangle. \end{aligned}$$

□

Since each term of Eq. (6.1) behaves differently and is controlled via different techniques depending on the result to be proven, we name each of these terms:

$$L_\tau := \langle \mu^\star - \mu_\tau^\star, \varepsilon \rangle, \quad Q_\tau := \|\Pi_\tau \varepsilon\|^2 \quad \text{and} \quad A_\tau := \|\mu^\star - \mu_\tau^\star\|^2. \quad (6.2)$$

It should be clear that L stands for “linear”, Q stands for “quadratic” and A stands for “approximation error”. We also define

$$\psi_\tau := 2L_\tau - Q_\tau + A_\tau. \quad (6.3)$$

Therefore a reformulation of Lemma 6.1 is

$$n\widehat{\mathcal{R}}_n(\tau) = \psi_\tau + \|\varepsilon\|^2.$$

Notice that $L_{\tau^\star} = A_{\tau^\star} = 0$ and $Q_{\tau^\star} \geq 0$, hence $\psi_{\tau^\star} \leq 0$. Also note that ψ , L and Q are random quantities depending on ε .

6.2. Deterministic bounds

In this section, we provide some deterministic bounds that are used in the proofs of Theorems 3.1 and 3.2.

6.2.1. Approximation error A_τ

We begin by the following result, which is the reason for the $\underline{\Lambda}_{\tau^\star} \underline{\Delta}^2$ term in Theorem 3.1.

Lemma 6.2. *Let $\tau \in \mathcal{T}_n$ be a segmentation such that $D := D_\tau < D_{\tau^\star}$. Then*

$$\frac{1}{n} A_\tau = \frac{1}{n} \|\mu^\star - \mu_\tau^\star\|^2 \geq \frac{1}{2} \underline{\Lambda}_{\tau^\star} \underline{\Delta}^2. \quad (6.4)$$

The proof of Lemma 6.2 can be found in Section B.3.2.

Remark 6.1. Lemma 6.2 is tight. Indeed, consider the simple case $D_\tau = 1$ and $D_{\tau^*} = 2$. Assume that $n = 2m$ is an even number, and let $\tau_1^* = m$. It follows from definitions (3.1) and (3.2) that, in this case,

$$\underline{\Delta} = \|\mu_1^* - \mu_n^*\|_{\mathcal{H}} \quad \text{and} \quad \underline{\Lambda}_{\tau^*} = \frac{1}{2}.$$

According to Eq. (2.5), $(\mu_\tau^*)_i = \frac{1}{2}(\mu_1^* + \mu_n^*)$, which yields

$$\frac{1}{n}A_\tau = \frac{1}{4}\|\mu_1^* - \mu_n^*\|_{\mathcal{H}}^2 = \frac{1}{2}\underline{\Lambda}_{\tau^*}\underline{\Delta}^2.$$

Thus, in this particular class of examples, equality holds in (6.4).

We next state an analogous result, valid for any $\tau \in \mathcal{T}_n$, which plays a key role in the proofs of Theorems 3.1 and 3.2.

Lemma 6.3. *For any $\tau \in \mathcal{T}_n$,*

$$\frac{1}{n}A_\tau \geq \frac{1}{2} \min \left\{ \underline{\Lambda}_{\tau^*}, \frac{1}{n}d_\infty^{(1)}(\tau^*, \tau) \right\} \underline{\Delta}^2. \quad (6.5)$$

Lemma 6.3 is proved in Section B.4.

6.2.2. Linear term L_τ and quadratic term Q_τ

The proof of Theorem 3.2 relies on some deterministic bounds on L_τ and Q_τ . We start with a preliminary lemma.

Lemma 6.4. *For any $\varepsilon_1, \dots, \varepsilon_n \in \mathcal{H}$,*

$$\frac{1}{2} \max_{1 \leq a < b \leq n} \left\| \sum_{j=a}^b \varepsilon_j \right\|_{\mathcal{H}} \leq \max_{1 \leq k \leq n} \left\| \sum_{j=1}^k \varepsilon_j \right\|_{\mathcal{H}} =: M_n. \quad (6.6)$$

Proof. For every $a < b$, we have:

$$\left\| \sum_{j=a}^b \varepsilon_j \right\|_{\mathcal{H}} = \left\| \sum_{j=1}^b \varepsilon_j - \sum_{j=1}^{a-1} \varepsilon_j \right\|_{\mathcal{H}} \leq \left\| \sum_{j=1}^b \varepsilon_j \right\|_{\mathcal{H}} + \left\| \sum_{j=1}^{a-1} \varepsilon_j \right\|_{\mathcal{H}} \leq 2M_n.$$

□

The following result is a deterministic bound on Q_τ in terms of M_n .

Lemma 6.5. *Let $\tau \in \mathcal{T}_n$ be a segmentation. Then*

$$Q_\tau \leq \frac{4D_\tau M_n^2}{n\underline{\Lambda}_\tau}.$$

Proof. By Eq. (2.5),

$$\begin{aligned}
Q_\tau &= \sum_{\ell=1}^{D_\tau} \frac{1}{|\tau_\ell - \tau_{\ell-1}|} \left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} \varepsilon_j \right\|_{\mathcal{H}}^2 \\
&\leq D_\tau \max_{1 \leq \ell \leq D_\tau} \left\{ \frac{1}{|\tau_\ell - \tau_{\ell-1}|} \left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} \varepsilon_j \right\|_{\mathcal{H}}^2 \right\} \\
&\leq \frac{D_\tau}{n\bar{\Delta}_\tau} \max_{1 \leq \ell \leq D_\tau} \left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} \varepsilon_j \right\|_{\mathcal{H}}^2 \leq \frac{4D_\tau}{n\bar{\Delta}_\tau} M_n^2,
\end{aligned}$$

where we used Lemma 6.4 for the last inequality. \square

The following result is a deterministic bound on L_τ .

Lemma 6.6. *For any $\tau \in \mathcal{T}_n$,*

$$|L_\tau| \leq 6D_{\tau^*} \max\{D_{\tau^*}, D_\tau\} \bar{\Delta} M_n.$$

Lemma 6.6 is proved in Section B.5.

6.3. Concentration

In this subsection, we present concentration results on Q_τ , L_τ , and deviation bounds for M_n — which will imply deviation bounds on Q_τ and L_τ by Lemmas 6.5 and 6.6). For any $j \in \{1, \dots, n\}$, $\tau \in \mathcal{T}_n$ and $\ell \in \{1, \dots, D_\tau\}$, we define

$$v_j := \mathbb{E} \left[\|\varepsilon_j\|_{\mathcal{H}}^2 \right] \quad v_{\tau, \ell} := \frac{1}{\tau_\ell - \tau_{\ell-1}} \sum_{j=\tau_{\ell-1}+1}^{\tau_\ell} v_j \quad \text{and} \quad v_\tau := \sum_{\ell=1}^{D_\tau} v_{\tau, \ell}.$$

Concentration under Assumption 1 The first result takes care of the linear term L_τ when Assumption 1 is satisfied.

Lemma 6.7 (Prop. 3 of Arlot, Celisse and Harchaoui [4]). *Suppose that Assumption 1 holds true. Then for any $x > 0$, with probability at least $1 - 2e^{-x}$, for any $\theta > 0$,*

$$|L_\tau| \leq \theta A_\tau + \left(\frac{4}{3} + \frac{1}{2\theta} \right) M^2 x.$$

The next result deals with the quadratic term Q_τ when Assumption 1 is satisfied.

Lemma 6.8 (Prop. 1 of Arlot, Celisse and Harchaoui [4]). *Suppose that Assumption 1 holds true. Then for any $x > 0$, with probability at least $1 - e^{-x}$,*

$$Q_\tau - v_\tau \leq \left(x + 2\sqrt{2xD_\tau} \right) \frac{14M^2}{3}.$$

We merge Lemmas 6.7 and 6.8 for convenience.

Lemma 6.9. *Suppose that Assumption 1 holds true. Take any $\lambda > 1$ and $\tau \in \mathcal{T}_n$ be a segmentation. Then, there exists an event $\Omega_{\tau,\lambda}^{(0)}$ of probability greater than $1 - 3e^{-\lambda D_\tau}$ on which:*

$$\psi_\tau \geq \frac{1}{3}A_\tau - \frac{74}{3}\lambda D_\tau M^2.$$

Proof. According to Lemma 6.7 with $\theta = 1/3$ and $x = \lambda D_\tau$, there exists an event $\Omega_{\tau,\lambda}^{(1)}$ on which $|L_\tau| \leq \frac{1}{3}A_\tau + \frac{17}{6}\lambda D_\tau M^2$, with $\mathbb{P}(\Omega_{\tau,\lambda}^{(1)}) \geq 1 - 2e^{-\lambda D_\tau}$. Lemma 6.8 with $x = \lambda D_\tau$ gives $\Omega_{\tau,\lambda}^{(2)}$ on which $Q_\tau - v_\tau \leq \frac{14}{3}(\lambda + 2\sqrt{2\lambda})D_\tau M^2$, with $\mathbb{P}(\Omega_{\tau,\lambda}^{(2)}) \geq 1 - e^{-\lambda D_\tau}$. Then, $\Omega_{\tau,\lambda}^{(0)} := \Omega_{\tau,\lambda}^{(1)} \cap \Omega_{\tau,\lambda}^{(2)}$ has a probability larger than $1 - 3e^{-\lambda D_\tau}$ by the union bound. Since for any $1 \leq \ell \leq D_\tau$, $v_{\tau,\ell} \leq M^2$, we have $v_\tau = \sum_{\ell=1}^{D_\tau} v_{\tau,\ell} \leq D_\tau M^2$. Hence, by definition (6.3) of ψ_τ and using that $\lambda \geq 1$, on the event $\Omega_{\tau,\lambda}^{(0)}$, we have:

$$\begin{aligned} \psi_\tau &\geq \frac{1}{3}A_\tau - \left(\frac{31}{3}\lambda + \frac{28}{3}\sqrt{2\sqrt{\lambda}} + 1 \right) D_\tau M^2 \\ &\geq \frac{1}{3}A_\tau - \lambda \left(\frac{31}{3} + \frac{28}{3}\sqrt{2} + 1 \right) D_\tau M^2. \end{aligned}$$

□

Remark 6.2. It is also possible to obtain an upper bound for ψ_τ : by Lemma 6.7, for every $\lambda \geq 0$, on the event $\Omega_{\tau,\lambda}^{(2)} \subset \Omega_{\tau,\lambda}^{(0)}$,

$$\psi_\tau \leq \frac{5}{3}A_\tau + \frac{17}{3}\lambda D_\tau M^2.$$

However, we do not need this result thereafter.

Concentration under Assumption 2 Lemma 6.5 and 6.6 directly translate upper bounds on M_n into controls of L_τ and Q_τ . Under Assumption 2, this is achieved via the following lemma, a Kolmogorov-like inequality for the noise in the RKHS. This result is a straightforward generalization of the inequality obtained by Kolmogorov [31] into the Hilbert setting. A more precise result (for real random variables only) can be found in [26], of which we follow the proof. The scheme of Hájek and Rényi [26] adapts well in our setting even though we do not need the full result.

Lemma 6.10. *If Assumption 2 holds true, then, for any $x > 0$,*

$$\mathbb{P}(M_n \geq x) \leq \frac{1}{x^2} \sum_{j=1}^n v_j. \quad (6.7)$$

We prove Lemma 6.10 in Section B.6.

Remark 6.3. We can reformulate Lemma 6.10 as follows. For any $y > 0$, there exists an event of probability at least $1 - y^{-2}$ on which $M_n < y\sqrt{\sum_{i=j}^n v_j} \leq y\sqrt{nV}$. Equivalently, for any $z \geq 0$, there exists an event of probability at least $1 - e^{-z}$ such that $M_n < e^{z/2} \sqrt{\sum_{i=j}^n v_j} \leq e^{z/2} \sqrt{nV}$.

6.4. Proof of Theorem 3.1

We follow the strategy described at the beginning of Section 6.

Definition of Ω Let us define $\Omega := \bigcap_{\tau \in \mathcal{T}_n} \Omega_{\tau, \lambda}^{(0)}$ with $\lambda = y + \log n + 1 > 1$, where we recall that $\Omega_{\tau, \lambda}^{(0)}$ is defined in Lemma 6.9. By the union bound, and since the $\Omega_{\tau, \lambda}^{(0)}$ have probability greater than $1 - 3e^{-\lambda D_\tau}$,

$$\mathbb{P}(\Omega) \geq 1 - 3 \sum_{\tau \in \mathcal{T}_n} e^{-\lambda D_\tau}.$$

The inequality $\mathbb{P}(\Omega) \geq 1 - e^{-y}$ follows since

$$\begin{aligned} \sum_{\tau \in \mathcal{T}_n} e^{-\lambda D_\tau} &= \sum_{d=1}^n \binom{n-1}{d-1} e^{-\lambda d} = e^{-\lambda} (1 + e^{-\lambda})^{n-1} \\ &\leq e^{-\lambda} \exp((n-1)e^{-\lambda}) \\ &= \frac{e^{-y}}{ne} \exp\left(\frac{n-1}{n} e^{-1-y}\right) \\ &\leq e^{-y} \frac{\exp(e^{-1})}{ne} \leq 0.27 e^{-y}, \end{aligned}$$

where the last inequality uses that $n \geq 2$. From now on we work exclusively on Ω .

Key argument We now make the simple (but crucial) observation that $\text{crit}(\tau^*) \geq \text{crit}(\hat{\tau})$, hence

$$n \text{pen}(\hat{\tau}) + \psi_{\hat{\tau}} \leq n \text{pen}(\tau^*) + \psi_{\tau^*} \leq n \text{pen}(\tau^*) = CD_{\tau^*} M^2.$$

Since we work on Ω , by definition of $\Omega_{\tau, \lambda}^{(0)}$ in Lemma 6.9, for any $\tau \in \mathcal{T}_n$, we have:

$$\psi_\tau \geq \frac{1}{3} A_\tau - \frac{74}{3} \lambda D_\tau M^2.$$

Therefore, we get:

$$CD_{\tau^*} M^2 \geq \frac{1}{3} A_{\hat{\tau}} + \left(C - \frac{74}{3} \lambda\right) D_{\hat{\tau}} M^2. \quad (6.8)$$

Proof that $D_{\hat{\tau}} \leq D_{\tau^*}$ Since $C > 74\lambda/3$ (by the lower bound in assumption (3.3)), $M^2 > 0$ and $A_{\hat{\tau}} \geq 0$, Eq. (6.8) implies that

$$D_{\hat{\tau}} \leq \frac{C}{C - \frac{74}{3}\lambda} D_{\tau^*}.$$

The lower bound in assumption (3.3) ensures that

$$\frac{C}{C - \frac{74}{3}\lambda} < \frac{D_{\tau^*} + 1}{D_{\tau^*}}$$

hence $D_{\hat{\tau}} \leq D_{\tau^*}$ on Ω .

Proof that $D_{\hat{\tau}} \geq D_{\tau^*}$ Since $C > 74\lambda/3$ (by the lower bound in assumption (3.3)), Eq. (6.8) implies that $A_{\hat{\tau}} \leq 3CD_{\tau^*}M^2$. A direct consequence of (3.3) is that $A_{\hat{\tau}} < \frac{1}{2}n\underline{\Lambda}_{\tau^*}\underline{\Delta}^2$, hence $D_{\hat{\tau}} \geq D_{\tau^*}$ by Lemma 6.2.

Loss between $\hat{\tau}$ and τ^* We have proved that $D_{\hat{\tau}} = D_{\tau^*}$ on Ω , therefore, Eq. (6.8) can be rewritten

$$A_{\hat{\tau}} \leq 74\lambda D_{\tau^*} M^2. \quad (6.9)$$

By Lemma 6.3 and the definition of λ , we get

$$\min \left\{ \underline{\Lambda}_{\tau^*}, \frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}) \right\} \leq \frac{148D_{\tau^*}M^2}{\underline{\Delta}^2} \cdot \frac{y + \log n + 1}{n} = v_1(y). \quad (6.10)$$

Remark that assumption (3.3) implies that

$$\frac{\underline{\Delta}^2}{M^2} \frac{\underline{\Lambda}_{\tau^*}}{6D_{\tau^*}} n > \frac{74}{3} (D_{\tau^*} + 1)(y + \log n + 1)$$

hence

$$\underline{\Lambda}_{\tau^*} > (D_{\tau^*} + 1) \frac{148D_{\tau^*}M^2}{\underline{\Delta}^2} \cdot \frac{y + \log n + 1}{n} > v_1(y).$$

Therefore, Eq. (6.10) can be simplified into

$$\frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}) \leq v_1(y). \quad \square$$

Remark 6.4. The proof of Theorem 3.1 generalizes to $\hat{\tau}$ defined by

$$\hat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n / \underline{\Lambda}_{\tau} \geq \delta_n} \{ \text{crit}(\tau) \}$$

instead of (2.2), for any $\delta_n \geq 0$ such that $\underline{\Lambda}_{\tau^*} \geq \delta_n$. Indeed, this assumption allows to write $\text{crit}(\tau^*) \geq \text{crit}(\hat{\tau})$ in the key argument, and the rest of the proof can stay unchanged (with the same event Ω). More generally, any constraint can be added in the argmin defining $\hat{\tau}$, provided that τ^* satisfies this constraint.

6.5. Proof of Theorem 3.2

We follow the strategy described at the beginning of Section 6. Throughout the proof, we write $\hat{\tau}_2$ as a shortcut for $\hat{\tau}(D_{\tau^*}, \delta_n)$.

Key argument By definition (3.6) of $\hat{\tau}_2 = \hat{\tau}(D_{\tau^*}, \delta_n)$, since we assume $\underline{\Lambda}_{\tau^*} \geq \delta_n$,

$$\hat{\mathcal{R}}_n(\tau^*) \geq \hat{\mathcal{R}}_n(\hat{\tau}_2)$$

hence

$$0 \geq \psi_{\tau^*} \geq \psi_{\hat{\tau}_2} = A_{\hat{\tau}_2} + 2L_{\hat{\tau}_2} - Q_{\hat{\tau}_2}.$$

By Lemma 6.5, Lemma 6.6 and the facts that $D_{\hat{\tau}_2} = D_{\tau^*}$ and $\underline{\Lambda}_{\hat{\tau}_2} \geq \delta_n$, we get

$$0 \geq \psi_{\hat{\tau}_2} \geq A_{\hat{\tau}_2} - 12D_{\tau^*}^2 \bar{\Delta} M_n - \frac{4D_{\tau^*} M_n^2}{n\delta_n}$$

hence, using Lemma 6.3,

$$\min \left\{ \underline{\Lambda}_{\tau^*}, \frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}_2) \right\} \leq \frac{24D_{\tau^*}^2 \bar{\Delta} M_n}{\underline{\Delta}^2 n} + \frac{8D_{\tau^*}}{\underline{\Delta}^2} \frac{M_n^2}{n^2 \delta_n}. \quad (6.11)$$

Definition of Ω_2 We define

$$\Omega_2 := \{M_n \leq y\sqrt{nV}\}.$$

By Lemma 6.10, under Assumption 2, $\mathbb{P}(\Omega_2) \geq 1 - y^{-2}$.

Conclusion By definition of Ω_2 , Eq. (6.11) implies that on Ω_2 :

$$\min \left\{ \underline{\Lambda}_{\tau^*}, \frac{1}{n} d_{\infty}^{(1)}(\tau^*, \hat{\tau}_2) \right\} \leq 24(D_{\tau^*})^2 \frac{\bar{\Delta}\sqrt{V}}{\underline{\Delta}^2} \frac{y}{\sqrt{n}} + 8D_{\tau^*} \frac{V}{\underline{\Delta}^2} \frac{y^2}{n\delta_n} = v_2(y, \delta_n).$$

Since we assume $v_2(y, \delta_n) < \underline{\Lambda}_{\tau^*}$, the result follows. \square

Appendix A: Additional notation

In this appendix are collected a large part of the technical details of the proofs that precede. Some additional notation used solely in the appendix are introduced below.

We denote by $\lambda_1^*, \dots, \lambda_{D_{\tau^*}}^*$ the segments of τ^* , that is,

$$\lambda_i^* = \{\tau_{i-1}^* + 1, \dots, \tau_i^*\}.$$

For any segment λ of $\tau \in \mathcal{T}_n$, we denote by μ_{λ}^* the value of μ_{τ}^* on λ , which does not depend on τ and is given by equation (2.5):

$$\mu_{\lambda}^* = \frac{1}{|\lambda|} \sum_{j \in \lambda} \mu_j^*. \quad (\text{A.1})$$

Appendix B: Proofs

B.1. Proof of Lemma 3.1

Proof of (i) We set $D^i := D_{\tau^i}$ for $i \in \{1, 2\}$. Let us show first that $d_\infty^{(2)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$. Take any $i \in \{1, \dots, D^1 - 1\}$, by the definition of $\underline{\Lambda}_{\tau^1}$,

$$|\tau_i^1 - \tau_{D^2}^2| = |\tau_i^1 - n| \geq n\underline{\Lambda}_{\tau^1} > n\underline{\Lambda}_{\tau^1}/2 \geq n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\}/2,$$

which is greater than $d_\infty^{(1)}(\tau^1, \tau^2)$ by assumption. In the same fashion we can prove that $|\tau_i^1 - \tau_0^2| > d_\infty^{(1)}(\tau^1, \tau^2)$. Hence, for any $i \in \{1, \dots, D^1 - 1\}$,

$$\min_{0 \leq j \leq D^2} |\tau_i^1 - \tau_j^2| = \min_{1 \leq j \leq D^2 - 1} |\tau_i^1 - \tau_j^2|,$$

which proves that $d_\infty^{(2)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$.

Next, we prove that $D^1 = D^2$ and $d_\infty^{(3)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$. Define $\phi : \{1, \dots, D^1 - 1\} \rightarrow \{1, \dots, D^2 - 1\}$ such that

$$\{\phi(i)\} = \arg \min_{1 \leq j \leq D^2 - 1} |\tau_i^1 - \tau_j^2|$$

for all $i \in \{1, \dots, D^1 - 1\}$. This mapping is well-defined: indeed, suppose that $j, k \in \{1, \dots, D^2 - 1\}$ both realize the minimum for some $i \in \{1, \dots, D^1 - 1\}$.

Since we assumed $\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) < \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\}/2$,

$$|\tau_i^1 - \tau_j^2| = |\tau_i^1 - \tau_k^2| \leq d_\infty^{(1)}(\tau^1, \tau^2) < n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\}/2.$$

By the triangle inequality,

$$|\tau_j^2 - \tau_k^2| < n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\} \leq n\underline{\Lambda}_{\tau^2},$$

hence $j = k$. Next, we show that ϕ is increasing. Take $i, j \in \{1, \dots, D^1 - 1\}$ such that $i < j$. Recall that τ^k is increasing ($k = 1, 2$). Then

$$\begin{aligned} \tau_{\phi(i)}^2 - \tau_{\phi(j)}^2 &= \tau_{\phi(i)}^2 - \tau_i^1 + \tau_i^1 - \tau_j^1 + \tau_j^1 - \tau_{\phi(j)}^2 \\ &= \tau_{\phi(i)}^2 - \tau_i^1 - |\tau_i^1 - \tau_j^1| + \tau_j^1 - \tau_{\phi(j)}^2 \\ &\leq \left| \tau_{\phi(i)}^2 - \tau_i^1 \right| - |\tau_i^1 - \tau_j^1| + \left| \tau_j^1 - \tau_{\phi(j)}^2 \right| \\ &\leq 2d_\infty^{(1)}(\tau^1, \tau^2) - |\tau_i^1 - \tau_j^1| \\ &< n \min\{\underline{\Lambda}_{\tau^1}, \underline{\Lambda}_{\tau^2}\} - n\underline{\Lambda}_{\tau^1} \leq 0. \end{aligned}$$

Hence $\phi(i) < \phi(j)$, so ϕ is increasing. As a consequence, ϕ is injective and we get $D^1 \leq D^2$. The same argument, exchanging τ^1 and τ^2 , shows that $D^2 \leq D^1$. Therefore, $D^1 = D^2$ and ϕ is an increasing permutation of $\{1, \dots, D^1 - 1\}$, hence it is the identity. As a consequence, $d_\infty^{(3)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^1, \tau^2)$.

Finally, since $d_\infty^{(3)}$ is symmetric, $d_\infty^{(i)}(\tau^1, \tau^2) = d_\infty^{(i)}(\tau^1, \tau^2)$ for any $i \in \{1, 2, 3\}$.

Proof of (ii) Since $D_{\tau^1} = D_{\tau^2}$, we can set $D = D_{\tau^1} = D_{\tau^2}$. Next, define $\phi(i) := \arg \min_{1 \leq j \leq D-1} |\tau_i^1 - \tau_j^2|$ and $C_\phi(i) := |\phi(i)|$ for all $i \in \{1, \dots, D-1\}$. Clearly, $C_\phi(i) \geq 1$ for any i . Let us show that we actually have $C_\phi(i) = 1$.

Take i and j distinct elements of $\{1, \dots, D-1\}$, and suppose that $\phi(i) \cap \phi(j)$ is non-empty. Let k be any element of $\phi(i) \cap \phi(j)$. By the triangle inequality and the definition of $d_\infty^{(1)}$,

$$n\Lambda_{\tau^1} \leq |\tau_i^1 - \tau_j^1| \leq |\tau_i^1 - \tau_k^2| + |\tau_k^2 - \tau_j^1| \leq 2d_\infty^{(1)}(\tau^1, \tau^2) < n\Lambda_{\tau^1}.$$

Hence, the $\phi(i)$ are disjoint and we can write $\sum_{i=1}^{D-1} C_\phi(i) = D-1$, which clearly implies that $C_\phi(i) = 1$.

From now on, we identify $\phi(i)$ with its unique element. Let us show that ϕ is increasing similarly to what we have done for proving (i). Take $i, j \in \{1, \dots, D-1\}$ such that $i < j$. We showed that

$$\tau_{\phi(i)}^2 - \tau_{\phi(j)}^2 \leq 2d_\infty^{(1)}(\tau^1, \tau^2) - |\tau_i^1 - \tau_j^1|,$$

thus according to the definition of Λ_{τ^1} , and our assumption,

$$\tau_{\phi(i)}^2 - \tau_{\phi(j)}^2 < n\Lambda_{\tau^1} - n\Lambda_{\tau^1} \leq 0.$$

Hence $\phi(i) < \phi(j)$: ϕ is increasing. As a consequence,

$$d_\infty^{(1)}(\tau^1, \tau^2) = d_\infty^{(1)}(\tau^2, \tau^1) = d_H^{(1)}(\tau^1, \tau^2).$$

□

B.2. The Frobenius loss

B.2.1. A formula for d_F^2

We start by proving a general formula for d_F , which is stated by Lajugie, Arlot and Bach [35], we prove it here for completeness:

$$\forall \tau^1, \tau^2 \in \mathcal{T}_n, \quad d_F(\tau^1, \tau^2)^2 = D_{\tau^1} + D_{\tau^2} - 2 \sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|^2}{|\lambda_k^1| \times |\lambda_\ell^2|}. \quad (\text{B.1})$$

Indeed, by definition, we have

$$d_F(\tau^1, \tau^2)^2 = \text{Tr}((\Pi_{\tau^1} - \Pi_{\tau^2})^2) = \underbrace{\text{Tr}(\Pi_{\tau^1})}_{=D_{\tau^1}} + \underbrace{\text{Tr}(\Pi_{\tau^2})}_{=D_{\tau^2}} - 2 \text{Tr}(\Pi_{\tau^1} \Pi_{\tau^2})$$

$$\text{and } \text{Tr}(\Pi_{\tau^1} \Pi_{\tau^2}) = \sum_{i=1}^n \sum_{j=1}^n \frac{\mathbf{1}_{\{\lambda_1(i)=\lambda_1(j) \text{ and } \lambda_2(i)=\lambda_2(j)\}}}{|\lambda_1(i)| |\lambda_2(i)|} = \sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|^2}{|\lambda_k^1| \times |\lambda_\ell^2|},$$

where we denoted by $\lambda_k(i)$ the segment of τ^k to which $i \in \{1, \dots, n\}$ belongs.

B.2.2. Proof of Eq. equation (3.4)

Eq. equation (3.4) is stated by Lajugie, Arlot and Bach [35]. The upper bound is a straightforward consequence of Eq. (B.1). We prove the lower bound here for completeness. We remark that

$$\sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|^2}{|\lambda_k^1| \times |\lambda_\ell^2|} \leq \sum_{k=1}^{D_{\tau^1}} \sum_{\ell=1}^{D_{\tau^2}} \frac{|\lambda_k^1 \cap \lambda_\ell^2|}{|\lambda_k^1|} = D_{\tau^1},$$

hence Eq. (B.1) shows that

$$d_F(\tau^1, \tau^2)^2 \geq D_{\tau^2} - D_{\tau^1}.$$

The lower bound follows since τ^1 and τ^2 play symmetric roles. \square

B.2.3. Proof of Proposition 3.1

Throughout the proof, we write $D = D_{\tau^1} = D_{\tau^2}$, $\epsilon = n^{-1}d_\infty^{(1)}(\tau^1, \tau^2)$ and we denote by $(\lambda_k^1)_{1 \leq k \leq D}$ and $(\lambda_k^2)_{1 \leq k \leq D}$ the segments of τ^1 and τ^2 , respectively.

Preliminary remark Since we assume that $D_{\tau^1} = D_{\tau^2}$ and $\frac{1}{n}d_\infty^{(1)}(\tau^1, \tau^2) = \epsilon < \underline{\Lambda}_{\tau^1}/2$, point (ii) in Lemma 3.1 shows that $d_\infty^{(1)}(\tau^1, \tau^2) = d_H^{(1)}(\tau^1, \tau^2) = d_\infty^{(3)}(\tau^1, \tau^2)$. In other words, for every $k \in \{1, \dots, D-1\}$, we have $|\tau_k^1 - \tau_k^2| \leq n\epsilon$, and some $k_0 \in \{1, \dots, D-1\}$ exists such that $|\tau_{k_0}^1 - \tau_{k_0}^2| = n\epsilon$. As a consequence, for every $k \in \{1, \dots, D-1\}$,

$$||\lambda_k^1| - |\lambda_k^2|| \leq 2n\epsilon \quad \text{and} \quad |\lambda_k^1 \cap \lambda_k^2| \geq |\lambda_k^1| - 2n\epsilon. \quad (\text{B.2})$$

Upper bound for $d_F(\tau^1, \tau^2)^2$ We focus on the sum appearing in the right-hand side of Eq. (B.1). Using Eq. (B.2), we get:

$$\begin{aligned} \sum_{k=1}^D \sum_{\ell=1}^D \frac{|\lambda_k^1 \cap \lambda_\ell^2|^2}{|\lambda_k^1| \times |\lambda_\ell^2|} &\geq \sum_{k=1}^D \frac{|\lambda_k^1 \cap \lambda_k^2|^2}{|\lambda_k^1| \times |\lambda_k^2|} \\ &\geq \sum_{k=1}^D \left[\frac{(|\lambda_k^1| - 2n\epsilon)^2}{|\lambda_k^1| \times (|\lambda_k^1| + 2n\epsilon)} \right] = \sum_{k=1}^D \frac{\left(1 - \frac{2n\epsilon}{|\lambda_k^1|}\right)^2}{1 + \frac{2n\epsilon}{|\lambda_k^1|}} \\ &\geq \sum_{k=1}^D \left(1 - \frac{6n\epsilon}{|\lambda_k^1|}\right) \geq D - \frac{6\epsilon D}{\underline{\Lambda}_{\tau^1}}, \end{aligned}$$

since for any $x \geq 0$, $\frac{(1-x)^2}{1+x} \geq 1 - 3x$. The upper bound follows, using Eq. (B.1).

Lower bound for $d_F(\tau^1, \tau^2)^2$ As shown in the preliminary remark, there exists some $k_0 \in \{1, \dots, D-1\}$ such that $|\tau_{k_0}^1 - \tau_{k_0}^2| = n\epsilon$. First consider the case where $\tau_{k_0}^1 < \tau_{k_0}^2$. Then, by definition of d_F and Π_τ , we have:

$$\begin{aligned} d_F(\tau^1, \tau^2)^2 &:= \sum_{1 \leq i, j \leq n} (\Pi_{\tau^1} - \Pi_{\tau^2})_{i,j}^2 \\ &\geq \sum_{i \in \lambda_{k_0+1}^1 \cap \lambda_{k_0}^2} \sum_{j \in \lambda_{k_0+1}^1 \cap \lambda_{k_0+1}^2} \frac{1}{|\lambda_{k_0+1}^1|^2} \\ &\quad + \sum_{i \in \lambda_{k_0+1}^1 \cap \lambda_{k_0+1}^2} \sum_{j \in \lambda_{k_0+1}^1 \cap \lambda_{k_0}^2} \frac{1}{|\lambda_{k_0+1}^1|^2} \\ &= \frac{2|\lambda_{k_0+1}^1 \cap \lambda_{k_0}^2| \cdot |\lambda_{k_0+1}^1 \cap \lambda_{k_0+1}^2|}{|\lambda_{k_0+1}^1|^2}. \end{aligned}$$

Now, remark that $|\lambda_{k_0+1}^1 \cap \lambda_{k_0}^2| = n\epsilon$, by the preliminary remark and our assumption $\tau_{k_0}^2 > \tau_{k_0}^1$. Using also Eq. (B.2), we get:

$$d_F(\tau^1, \tau^2)^2 \geq \frac{2n\epsilon(|\lambda_{k_0+1}^1| - 2n\epsilon)}{|\lambda_{k_0+1}^1|^2} \geq \frac{2n\epsilon}{3\bar{\Lambda}_{\tau^1}},$$

since $|\lambda_{k_0+1}^1| - 2n\epsilon \geq |\lambda_{k_0+1}^1|/3$ and $|\lambda_{k_0+1}^1| \leq \bar{\Lambda}_{\tau^1}$. When $\tau_{k_0}^1 > \tau_{k_0}^2$, we apply the same reasoning, restricting the sum over i, j in the definition of d_F to $i \in \lambda_{k_0}^1 \cap \lambda_{k_0}^2$ and $j \in \lambda_{k_0}^1 \cap \lambda_{k_0+1}^2$ (plus its symmetric). We obtain the same lower bound, which concludes the proof. \square

B.3. Lower bounds on the approximation error

This section provides the proofs of Lemmas 6.2 and 6.3.

B.3.1. Preliminary lemma

We start with a lemma useful in the two proofs.

Lemma B.1. *If a segment $\lambda \subset \{1, \dots, n\}$ intersects only two segments of τ^* , λ_i^* and λ_{i+1}^* , then we have:*

$$\sum_{j \in \lambda} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 = \frac{|\lambda \cap \lambda_i^*| \cdot |\lambda \cap \lambda_{i+1}^*|}{|\lambda \cap \lambda_i^*| + |\lambda \cap \lambda_{i+1}^*|} \left\| \mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^* \right\|_{\mathcal{H}}^2 \quad (\text{B.3})$$

$$\geq \left(\frac{|\lambda \cap \lambda_i^*|}{|\lambda_i^*|} \wedge \frac{|\lambda \cap \lambda_{i+1}^*|}{|\lambda_{i+1}^*|} \right) \cdot \frac{|\lambda_i^*| \cdot |\lambda_{i+1}^*|}{|\lambda_i^*| + |\lambda_{i+1}^*|} \cdot \left\| \mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^* \right\|_{\mathcal{H}}^2. \quad (\text{B.4})$$

Proof. We first prove Eq. (B.3). Since λ only intersects λ_i^* and λ_{i+1}^* , we have:

$$\begin{aligned} \sum_{j \in \lambda} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 &= \sum_{j \in \lambda \cap \lambda_i^*} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 + \sum_{j \in \lambda \cap \lambda_{i+1}^*} \|\mu_j^* - \mu_\lambda^*\|_{\mathcal{H}}^2 \\ &= |\lambda \cap \lambda_i^*| \cdot \|\mu_{\lambda_i^*}^* - \mu_\lambda^*\|_{\mathcal{H}}^2 + |\lambda \cap \lambda_{i+1}^*| \cdot \|\mu_{\lambda_{i+1}^*}^* - \mu_\lambda^*\|_{\mathcal{H}}^2. \end{aligned} \quad (\text{B.5})$$

Since μ_λ^* is given by Eq. (A.1), we obtain

$$\begin{aligned} \|\mu_{\lambda_i^*}^* - \mu_\lambda^*\|_{\mathcal{H}}^2 &= \left\| \frac{1}{|\lambda|} \sum_{j \in \lambda} (\mu_{\lambda_i^*}^* - \mu_j^*) \right\|_{\mathcal{H}}^2 = \left\| \frac{1}{|\lambda|} \sum_{j \in \lambda \cap \lambda_{i+1}^*} (\mu_{\lambda_i^*}^* - \mu_{\lambda_{i+1}^*}^*) \right\|_{\mathcal{H}}^2 \\ &= \frac{|\lambda \cap \lambda_{i+1}^*|^2}{|\lambda|^2} \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2. \end{aligned}$$

The same computation on $\lambda \cap \lambda_{i+1}^*$ yields

$$\|\mu_{\lambda_{i+1}^*}^* - \mu_\lambda^*\|_{\mathcal{H}}^2 = \frac{|\lambda \cap \lambda_i^*|^2}{|\lambda|^2} \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2.$$

Therefore, Eq. (B.5) and the fact that $|\lambda| = |\lambda \cap \lambda_i^*| + |\lambda \cap \lambda_{i+1}^*|$ yield Eq. (B.3).

Now, we remark that for any $a, b, c, d > 0$,

$$\frac{abcd}{ab + cd} = \frac{1}{\frac{ab}{\max(a, c)} + \frac{cd}{\max(a, c)}} \times \min(a, c) \times bd \geq \min(a, c) \frac{bd}{b + d}.$$

Taking $a = |\lambda \cap \lambda_i^*| / |\lambda_i^*|$, $b = |\lambda_i^*|$, $c = |\lambda \cap \lambda_{i+1}^*| / |\lambda_{i+1}^*|$ and $d = |\lambda_{i+1}^*|$, we get Eq. (B.4). \square

B.3.2. Proof of Lemma 6.2

In fact, we prove a slightly stronger statement. We show that, for any $n \geq 2$, for any $D_{\tau^*} \in \{2, \dots, n\}$, for any $D \in \{1, \dots, D_{\tau^*} - 1\}$ and any $\tau \in \mathcal{T}_n^D$,

$$\|\mu^* - \mu_\tau^*\|^2 \geq \min_{1 \leq i \leq D_{\tau^*} - 1} \left\{ \frac{|\lambda_i^*| \cdot |\lambda_{i+1}^*|}{|\lambda_i^*| + |\lambda_{i+1}^*|} \cdot \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2 \right\}. \quad (\text{B.6})$$

Then,

$$\|\mu^* - \mu_\tau^*\|^2 \geq \underline{\Gamma} \cdot \underline{\Delta}^2 \quad \text{where} \quad \underline{\Gamma} = \left(n \max_{1 \leq i \leq D_{\tau^*} - 1} \left\{ \frac{1}{|\lambda_i^*|} + \frac{1}{|\lambda_{i+1}^*|} \right\} \right)^{-1}.$$

Since we always have

$$\underline{\Delta}_{\tau^*} \geq \underline{\Gamma} \geq \frac{1}{2} \underline{\Delta}_{\tau^*},$$

Eq. equation (6.4) follows.

Proof of Eq. (B.6) by induction We show by strong induction on D_{τ^*} that, for any $D_{\tau^*} \geq 2$, for any $D \in \{1, \dots, D_{\tau^*} - 1\}$, any $n \geq D_{\tau^*}$ and any $\tau \in \mathcal{T}_n^D$, Eq. (B.6) holds true.

First, if $D_{\tau^*} = 2$, the result follows by Eq. (B.4) in Lemma B.1 since we then have $i = 1$ and

$$\frac{|\lambda \cap \lambda_1^*|}{|\lambda_1^*|} = \frac{|\lambda \cap \lambda_2^*|}{|\lambda_2^*|} = 1.$$

Suppose now that the result is proved for all $D_{\tau^*} \in \{2, \dots, p\}$ and consider a change-point problem (τ^*, μ^*) with $D_{\tau^*} = D_{\tau^*} = p + 1$ and $n \geq p + 1$. Let $D < p + 1$ and some segmentation $\tau \in \mathcal{T}_n^D$ be fixed. Then one of these two scenarios occurs: (i) there exists λ_i^* with $2 \leq i \leq D_{\tau^*} - 1$ that does not contain any change-point of τ , or (ii) $\lambda_2^*, \dots, \lambda_{D_{\tau^*}-1}^*$ all contain a change-point of τ .

Case (i) Suppose that there exists an inner segment λ_i^* of τ^* , $2 \leq i \leq D_{\tau^*} - 1$, that does not contain any change-point of τ (see Figure 7). Therefore, there exists $k \in \{1, \dots, D\}$ such that $\lambda_i^* \subsetneq \lambda_k$. By definition, there are $i - 1$ change-points of τ^* to the left of λ_i^* and $k - 1$ change-points of τ to the left of λ_i^* . Suppose that $k < i$. We define τ° as the segmentation obtained by adding τ_i^* to τ (see Figure 7). Then $\|\mu^* - \mu_{\tau^\circ}^*\|^2 \geq \|\mu^* - \mu_{\tau^*}^*\|^2$ because τ° is finer than τ . Reducing τ° to a segmentation $\tilde{\tau}^\circ$ of $\{1, 2, \dots, \tau_i^*\}$ in k segments and τ^* to a segmentation $\tilde{\tau}^*$ of $\{1, 2, \dots, \tau_i^*\}$ in i segments and defining $\tilde{\mu}^* = (\mu_1^*, \dots, \mu_{\tau_i^*}^*) \in \mathcal{H}^i$, we get back to a situation covered by the induction since $i \leq D_{\tau^*} - 1$ and $k < i$. So,

$$\begin{aligned} \|\tilde{\mu}^* - \tilde{\mu}_{\tilde{\tau}^\circ}^*\|^2 &\geq \inf_{1 \leq j \leq i-1} \left\{ \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \tilde{\mu}_{\lambda_{j+1}^*}^* - \tilde{\mu}_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right\} \\ &\geq \inf_{1 \leq j \leq D_{\tau^*}-1} \left\{ \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \mu_{\lambda_{j+1}^*}^* - \mu_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right\} \end{aligned}$$

and we get the result since $\|\mu^* - \mu_{\tau^\circ}^*\|^2 \geq \|\tilde{\mu}^* - \tilde{\mu}_{\tilde{\tau}^\circ}^*\|^2$. A symmetric reasoning can be applied if $k \geq i$, considering change-points to the right of λ_i^* and using that $D - k + 1 < D_{\tau^*} - i + 1$ since $D < D_{\tau^*}$.

| | | | |
|----------------------|-----|---------------|-----|
| $\tilde{\tau}^*$ | ... | | |
| τ^* | ... | λ_i^* | ... |
| τ | ... | λ_k | ... |
| τ° | ... | | ... |
| $\tilde{\tau}^\circ$ | ... | | |

FIG 7. Proof of Lemma 6.2, Case (i): λ_i^* is a segment of τ^* that is included in a segment of τ . The segmentation τ° is obtained by joining τ_i^* to the segmentation τ .

Case (ii) Suppose that each inner segment of τ^* contains a change-point of τ . Since there are $D_{\tau^*} - 2$ inner segments of τ^* and $D - 1 \leq D_{\tau^*} - 2$ change-points

of τ , there is at most (hence exactly) one change-point of τ in each inner segment of τ^* . Then $D = D_{\tau^*} - 1$ and we are in the situation depicted in Figure 8.

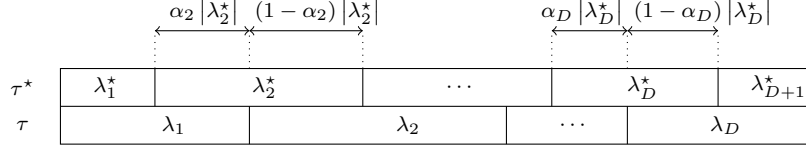


FIG 8. *Proof of Lemma 6.2, Case (ii): $D = D_{\tau^*} - 1$ and each inner segment of τ^* contains exactly one change-point of τ .*

We can use Eq. (B.4) in Lemma B.1 to lower bound the contribution of each $\lambda \in \tau$ to $\|\mu^* - \mu_\tau^*\|^2$. For $2 \leq i \leq D = D_{\tau^*} - 1$, define $\alpha_i := |\lambda_i^* \cap \lambda_{i-1}| / |\lambda_i^*|$. Then, we have

$$\begin{aligned}
\|\mu^* - \mu_\tau^*\|^2 &\geq (1 \wedge \alpha_2) \frac{|\lambda_1^*| \cdot |\lambda_2^*|}{|\lambda_1^*| + |\lambda_2^*|} \cdot \left\| \mu_{\lambda_2^*}^* - \mu_{\lambda_1^*}^* \right\|_{\mathcal{H}}^2 \\
&\quad + \sum_{j=2}^{D-1} \left([(1 - \alpha_j) \wedge \alpha_{j+1}] \cdot \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \mu_{\lambda_{j+1}^*}^* - \mu_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right) \\
&\quad + [(1 - \alpha_D) \wedge 1] \frac{|\lambda_D^*| \cdot |\lambda_{D+1}^*|}{|\lambda_D^*| + |\lambda_{D+1}^*|} \cdot \left\| \mu_{\lambda_{D+1}^*}^* - \mu_{\lambda_D^*}^* \right\|_{\mathcal{H}}^2 \\
&\geq [1 \wedge \alpha_2 + (1 - \alpha_2) \wedge \alpha_3 + \cdots + (1 - \alpha_{D-1}) \wedge \alpha_D + (1 - \alpha_D) \wedge 1] \\
&\quad \times \inf_{1 \leq j \leq D_{\tau^*} - 1} \left\{ \frac{|\lambda_j^*| \cdot |\lambda_{j+1}^*|}{|\lambda_j^*| + |\lambda_{j+1}^*|} \cdot \left\| \mu_{\lambda_{j+1}^*}^* - \mu_{\lambda_j^*}^* \right\|_{\mathcal{H}}^2 \right\}.
\end{aligned}$$

Since $\alpha_i \geq 0$ for any $2 \leq i \leq D_{\tau^*} - 1$, it is straightforward to show that

$$\alpha_2 + (1 - \alpha_2) \wedge \alpha_3 + \cdots + (1 - \alpha_D) \geq 1,$$

which concludes the proof. \square

B.4. Proof of Lemma 6.3

Let us define $\delta := \min\{n_{\underline{\Lambda}_{\tau^*}}, d_{\infty}^{(1)}(\tau^*, \tau)\}$. If $\delta = 0$, then Eq. equation (6.5) holds true. We assume from now on that $\delta > 0$.

Because $n_{\underline{\Lambda}_{\tau^*}} \geq \delta$, for any $1 \leq i \leq D_{\tau^*} - 1$, we can write $|\tau_{i+1}^* - \tau_i^*| \geq \delta$. On the other hand, because $d_{\infty}^{(1)}(\tau^*, \tau) \geq \delta$, there exists $i \in \{1, \dots, D_{\tau^*} - 1\}$ such that, for any $j \in \{1, \dots, D - 1\}$, $|\tau_i^* - \tau_j| \geq \delta$. Since $\delta \leq n_{\underline{\Lambda}_{\tau^*}}$, this also holds true for $j = 0$ and $j = D$. Let us define, as illustrated by Figure 9,

$$\lambda^\circ := \{\tau_i^* - \delta + 1, \dots, \tau_i^*, \tau_i^* + 1, \dots, \tau_i^* + \delta\}.$$

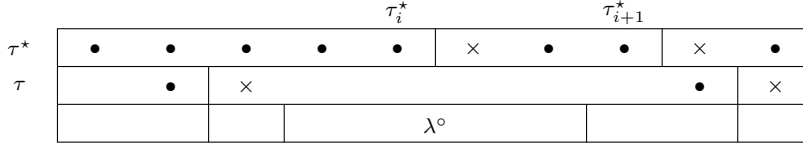


FIG 9. Construction of λ° in the proof of Lemma 6.3. In this case, $\delta = 2$ since $\underline{\Delta}_{\tau^*} = 2/10$ (the rightmost segment of τ^* is of size 2) and $d_\infty(\tau^*, \tau) = 3$ (achieved in τ_i^*).

Since λ° is included in a segment of τ ,

$$\|\mu^* - \mu_\tau^*\|^2 \geq \sum_{j \in \lambda^\circ} \|\mu_j^* - (\mu_\tau^*)_j\|_{\mathcal{H}}^2 \geq \sum_{j \in \lambda^\circ} \|\mu_j^* - \mu_{\lambda^\circ}^*\|_{\mathcal{H}}^2.$$

Because of the hypothesis we made, λ° only intersects λ_i^* and λ_{i+1}^* among the segments of τ^* , so Eq. (B.3) in Lemma B.1 shows that

$$\begin{aligned} \sum_{j \in \lambda^\circ} \|\mu_j^* - \mu_{\lambda^\circ}^*\|_{\mathcal{H}}^2 &= \frac{|\lambda^\circ \cap \lambda_i^*| \cdot |\lambda^\circ \cap \lambda_{i+1}^*|}{|\lambda^\circ \cap \lambda_i^*| + |\lambda^\circ \cap \lambda_{i+1}^*|} \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2 \\ &= \frac{\delta}{2} \|\mu_{\lambda_{i+1}^*}^* - \mu_{\lambda_i^*}^*\|_{\mathcal{H}}^2 \geq \frac{\delta}{2} \Delta^2, \end{aligned}$$

hence the result. \square

B.5. Proof of Lemma 6.6

In this proof, since τ is fixed, we denote by $\lambda_1, \dots, \lambda_D$ the segments of τ , that is, $\lambda_i = \{\tau_{i-1} + 1, \dots, \tau_i\}$.

First, notice that

$$L_\tau = \langle \mu^* - \mu_\tau^*, \varepsilon \rangle = \sum_{i=1}^{D_{\tau^*}} \left\langle \mu_{\lambda_i^*}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} - \sum_{i=1}^{D_\tau} \left\langle \mu_{\lambda_i}^*, \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}}. \quad (\text{B.7})$$

Now, if $D_\tau < D_{\tau^*}$ we *arbitrarily* define $\lambda_{D_\tau+1} = \dots = \lambda_{D_{\tau^*}} = \emptyset$, so that $\sum_{j \in \lambda_i} \varepsilon_j = 0$ for every $i \in \{D_\tau + 1, \dots, D_{\tau^*}\}$. Similarly, if $D_{\tau^*} < D_\tau$, we define $\lambda_{D_{\tau^*}+1}^* = \dots = \lambda_{D_\tau}^* = \emptyset$. We also define $\mu_\emptyset^* = \mu_n^*$ by convention. Then, defining $D^+ := \max\{D_{\tau^*}, D_\tau\}$, we can rewrite Eq. (B.7) as follows:

$$\begin{aligned} L_\tau &= \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i^*}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} - \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i}^*, \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i^*}^* - \mu_{\lambda_i}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} + \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i}^*, \sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}} \\ &= \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i^*}^* - \mu_{\lambda_i}^*, \sum_{j \in \lambda_i^*} \varepsilon_j \right\rangle_{\mathcal{H}} + \sum_{i=1}^{D^+} \left\langle \mu_{\lambda_i}^* - \mu_n^*, \sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right\rangle_{\mathcal{H}}, \end{aligned}$$

since

$$\sum_{i=1}^{D^+} \left(\sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right) = 0.$$

Then, by the triangle inequality and Cauchy-Schwarz inequality,

$$\begin{aligned} |L_\tau| &\leq \sum_{i=1}^{D^+} \left\| \mu_{\lambda_i^*}^* - \mu_{\lambda_i}^* \right\|_{\mathcal{H}} \left\| \sum_{j \in \lambda_i^*} \varepsilon_j \right\|_{\mathcal{H}} + \sum_{i=1}^{D^+} \left\| \mu_{\lambda_i}^* - \mu_n^* \right\|_{\mathcal{H}} \left\| \sum_{j \in \lambda_i^*} \varepsilon_j - \sum_{j \in \lambda_i} \varepsilon_j \right\|_{\mathcal{H}} \\ &\leq \text{diam conv} \{ \mu_j^* / j \in \{1, \dots, n\} \} \\ &\quad \times \left[\sum_{i=1}^{D^+} \left\| \sum_{j \in \lambda_i^*} \varepsilon_j \right\|_{\mathcal{H}} + \sum_{i=1}^{D^+} \left(\left\| \sum_{j \in \lambda_i^*} \varepsilon_j \right\|_{\mathcal{H}} + \left\| \sum_{j \in \lambda_i} \varepsilon_j \right\|_{\mathcal{H}} \right) \right] \\ &\leq 3D^+ \text{diam conv} \{ \mu_j^* / j \in \{1, \dots, n\} \} \times \sup_{1 \leq a < b \leq n} \left\| \sum_{j=a}^b \varepsilon_j \right\|_{\mathcal{H}} \end{aligned}$$

where we used that $\mu_\lambda^* \in \text{conv} \{ \mu_j^* / j \in \{1, \dots, n\} \}$ for any segment λ . Since the diameter of the convex hull of a finite set of points is equal to the diameter of the set, we have

$$\begin{aligned} \text{diam conv} \{ \mu_j^* / j \in \{1, \dots, n\} \} &= \text{diam} \{ \mu_j^* / j \in \{1, \dots, n\} \} \\ &\leq (D_{\tau^*} - 1) \overline{\Delta} < D_{\tau^*} \overline{\Delta}. \end{aligned}$$

Using also Lemma 6.4, we get the result. \square

B.6. Proof of Lemma 6.10

Let us put $\zeta := \|\varepsilon_1 + \dots + \varepsilon_n\|_{\mathcal{H}}^2$. Since for any $j \neq k$, $\mathbb{E} [\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}}] = 0$ (see Remark B.1), by definition of v_j ,

$$\mathbb{E} [\zeta] = \mathbb{E} \left[\|\varepsilon_1 + \dots + \varepsilon_n\|_{\mathcal{H}}^2 \right] = \sum_{j=1}^n v_j.$$

We recognize the right-hand side of equation (6.7) up to $1/x^2$. For any $r > 1$, let us denote by A_r the event

$$\forall 1 \leq s < r, \quad \|\varepsilon_1 + \dots + \varepsilon_s\|_{\mathcal{H}} < x \quad \text{and} \quad \|\varepsilon_1 + \dots + \varepsilon_r\|_{\mathcal{H}} \geq x,$$

and by A_1 the event $\|\varepsilon_1\|_{\mathcal{H}} \geq x$. These events are disjoint, thus we can write

$$\mathbb{P} \left(\max_{1 \leq k \leq n} \|\varepsilon_1 + \dots + \varepsilon_k\|_{\mathcal{H}} \geq x \right) = \mathbb{P} \left(\bigcup_{r=1}^n A_r \right) = \sum_{r=1}^n \mathbb{P}(A_r). \quad (\text{B.8})$$

The law of total expectation and the positiveness of ζ yield

$$\mathbb{E}[\zeta] \geq \sum_{r=1}^n \mathbb{E}[\zeta | A_r] \mathbb{P}(A_r) .$$

Finally, let $\ell \leq r < k$ be integers. Since ε_ℓ is independent from ε_k conditionally to $\sigma(\varepsilon_1, \dots, \varepsilon_r)$, ε_ℓ is independent from ε_k conditionally to A_r . Furthermore, ε_k is independent from A_r and

$$\mathbb{E}[\langle \varepsilon_k, \varepsilon_\ell \rangle_{\mathcal{H}} | A_r] = \langle \mathbb{E}[\varepsilon_k], \mathbb{E}[\varepsilon_\ell | A_r] \rangle_{\mathcal{H}} = 0 .$$

Because of this relation and the positivity of the (real) conditional expectation, for any integers $r \leq k \leq j$,

$$\mathbb{E}[\zeta | A_r] = \mathbb{E}[\|\varepsilon_1 + \dots + \varepsilon_n\|_{\mathcal{H}}^2 | A_r] \geq \mathbb{E}[\|\varepsilon_1 + \dots + \varepsilon_r\|_{\mathcal{H}}^2 | A_r] \geq x^2 .$$

Therefore, $\mathbb{E}[\zeta | A_r] \geq x^2$, which gives $\mathbb{E}[\zeta] \geq x^2 \sum \mathbb{P}(A_r)$. This concludes the proof, thanks to Eq. (B.8). \square

Remark B.1. The independence between ε_j and ε_k for $j \neq k$ yields $\mathbb{E}[\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}}] = 0$. Indeed, we dispose of a conditional expectation on \mathcal{H} [19, chapter 5], which satisfies the same properties than the conditional expectation with real random variables. Hence we can write

$$\begin{aligned} \mathbb{E}[\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}}] &= \mathbb{E}[\mathbb{E}[\langle \varepsilon_j, \varepsilon_k \rangle_{\mathcal{H}} | \varepsilon_k]] = \mathbb{E}[\langle \mathbb{E}[\varepsilon_j | \varepsilon_k], \varepsilon_k \rangle_{\mathcal{H}}] \\ &= \mathbb{E}[\langle \mathbb{E}[\varepsilon_j], \varepsilon_k \rangle_{\mathcal{H}}] = 0 . \end{aligned}$$

Note that the ε_j s expectation vanishes by hypothesis.

Acknowledgments

Damien Garreau PhD scholarship is financed by DGA / Inria. Sylvain Arlot is also member of the Select project-team of Inria Saclay. At the beginning of this work, Sylvain Arlot was financed by CNRS and member of the Sierra team in the Département d'Informatique de l'École normale supérieure (CNRS / ENS / Inria UMR 8548), 45 rue d'Ulm, 75005 Paris, France. This work was also partly done while Sylvain Arlot was supported by Institut des Hautes Études Scientifiques (IHES, Le Bois-Marie, 35, route de Chartres, 91440 Bures-Sur-Yvette, France). The authors thank Alain Celisse and Aymeric Dieuleveut for helpful discussions.

References

- [1] ABOU-ELAILAH, A., GOUET-BRUNET, V. and BLOCH, I. (2015). Detection of Abrupt Changes in Spatial Relationships in Video Sequences. In *International Conference on Pattern Recognition Applications and Methods* 89–106. Springer.

- [2] ABRAMOVICH, F., BENJAMINI, Y., DONOHO, D. L. and JOHNSTONE, I. M. (2006). Adapting to unknown sparsity by controlling the false discovery rate. *Annals of Statistics* **34** 584–653.
- [3] ARLOT, S. and CELISSE, A. (2011). Segmentation of the mean of heteroscedastic data via cross-validation. *Statistics and Computing* **21** 613–632.
- [4] ARLOT, S., CELISSE, A. and HARCHAOUI, Z. (2012). A kernel multiple change-point algorithm via model selection. *ArXiv e-prints*. Available at <https://arxiv.org/abs/1202.3878v2>.
- [5] ARONSAJN, N. (1950). Theory of reproducing kernels. *Transactions of the American mathematical society* 337–404.
- [6] BAI, J. and PERRON, P. (1998). Estimating and testing linear models with multiple structural changes. *Econometrica* 47–78.
- [7] BASSEVILLE, M. and NIKIFOROV, I. V. (1993). *Detection of abrupt changes: theory and application*. Prentice Hall Englewood Cliffs.
- [8] BAUDRY, J.-P., MAUGIS, C. and MICHEL, B. (2012). Slope heuristics: overview and implementation. *Statistics and Computing* **22** 455–470.
- [9] BELLMAN, R. (1961). On the approximation of curves by line segments using dynamic programming. *Communications of the ACM* **4** 284.
- [10] BIRGÉ, L. and MASSART, P. (2001). Gaussian model selection. *Journal of the European Mathematical Society* **3** 203–268.
- [11] BIRGÉ, L. and MASSART, P. (2007). Minimal penalties for Gaussian model selection. *Probability Theory and Related Fields* **138** 33–73.
- [12] BOYSEN, L., KEMPE, A., LIEBSCHER, V., MUNK, A. and WITTICH, O. (2009). Consistencies and rates of convergence of jump-penalized least squares estimators. *Annals of Statistics* **37** 157–183.
- [13] BRODSKY, B. E. and DARKHOVSKY, B. S. (2013). *Nonparametric methods in change point problems* **243**. Springer Science & Business Media.
- [14] BRUNEL, V.-E. (2014). Convex set detection. *ArXiv e-prints*. Available at <https://arxiv.org/abs/1404.6224>.
- [15] CARLSTEIN, E. (1988). Nonparametric change-point estimation. *Annals of Statistics* 188–197.
- [16] CELISSE, A., MAROT, G., PIERRE-JEAN, M. and RIGAILL, G. (2016). New efficient algorithms for multiple change-point detection with kernels. *Preprint*. Available at <https://hal.inria.fr/hal-01413230>.
- [17] COMTE, F. and ROZENHOLC, Y. (2004). A new algorithm for fixed design regression and denoising. *Annals of the Institute of Statistical Mathematics* **56** 449–473.
- [18] DESOBRY, F., DAVY, M. and DONCARLI, C. (2005). An online kernel change detection algorithm. *Signal Processing, IEEE Transactions on* **53** 2961–2974.
- [19] DIESTEL, J. and UHL, J. J. (1977). *Vector measures* **15**. American Mathematical Soc.
- [20] DIEULEVEUT, A. and BACH, F. (2016). Nonparametric stochastic approximation with large step-sizes. *Annals of Statistics* **44** 1363–1399.
- [21] FISHER, W. D. (1958). On grouping for maximum homogeneity. *Journal*

- of the American Statistical Association **53** 789–798.
- [22] FRYZLEWICZ, P. (2014). Wild binary segmentation for multiple change-point detection. *The Annals of Statistics* **42** 2243–2281.
 - [23] FUKUMIZU, K., GRETTON, A., SUN, X. and SCHÖLKOPF, B. (2008). Kernel Measures of Conditional Dependence. In *Advances in Neural Information Processing Systems 20* 489–496. Curran Associates, Inc.
 - [24] GRETTON, A., BORGWARDT, K. M., RASCH, M., SCHÖLKOPF, B. and SMOLA, A. J. (2006). A kernel method for the two-sample-problem. In *Advances in neural information processing systems* 513–520.
 - [25] GRETTON, A., SEJDINOVIC, D., STRATHMANN, H., BALAKRISHNAN, S., PONTIL, M., FUKUMIZU, K. and SRIPERUMBUDUR, B. K. (2012). Optimal kernel choice for large-scale two-sample tests. In *Advances in Neural Information Processing Systems* 1205–1213.
 - [26] HÁJEK, J. and RÉNYI, A. (1955). Generalization of an inequality of Kolmogorov. *Acta Mathematica Hungarica* **6** 281–283.
 - [27] HARCHAOUI, Z. and CAPPÉ, O. (2007). Retrospective multiple change-point estimation with kernels. In *IEEE Workshop on Statistical Signal Processing* 768–772.
 - [28] HARCHAOUI, Z., MOULINES, E. and BACH, F. R. (2009). Kernel change-point analysis. In *Advances in neural information processing systems* 609–616.
 - [29] HUBERT, L. and ARABIE, P. (1985). Comparing partitions. *Journal of classification* **2** 193–218.
 - [30] KIM, A. Y., MARZBAN, C., PERCIVAL, D. B. and STUETZLE, W. (2009). Using labeled data to evaluate change detectors in a multivariate streaming environment. *Signal Processing* **89** 2529–2536.
 - [31] KOLMOGOROV, A. N. (1928). Über die Summen durch den Zufall bestimmten unabhängigen Größen. *Mathematische Annalen* **99** 484–488.
 - [32] KOROSTELEV, A. P. (1988). On minimax estimation of a discontinuous signal. *Theory of Probability & Its Applications* **32** 727–730.
 - [33] KOROSTELEV, A. P. and TSYBAKOV, A. B. (2012). *Minimax theory of image reconstruction* **82**. Springer Science & Business Media.
 - [34] LAI, W. R., JOHNSON, M. D., KUCHERLAPATI, R. and PARK, P. J. (2005). Comparative analysis of algorithms for identifying amplifications and deletions in array CGH data. *Bioinformatics* **21** 3763–3770.
 - [35] LAJUGIE, R., ARLOT, S. and BACH, F. (2014). Large-margin metric learning for constrained partitioning problems. In *Proceedings of The 31st International Conference on Machine Learning* 297–305.
 - [36] LAVIELLE, M. (2005). Using penalized contrasts for the change-point problem. *Signal processing* **85** 1501–1510.
 - [37] LAVIELLE, M. and MOULINES, E. (2000). Least-squares Estimation of an Unknown Number of Shifts in a Time Series. *Journal of time series analysis* **21** 33–59.
 - [38] LAVIELLE, M. and TEYSSIERE, G. (2006). Detection of multiple change-points in multivariate time series. *Lithuanian Mathematical Journal* **46** 287–306.

- [39] LEBARBIER, É. (2005). Detecting multiple change-points in the mean of a Gaussian process by model selection. *Signal Proces.* **85** 717–736.
- [40] LEDOUX, M. and TALAGRAND, M. (2013). *Probability in Banach Spaces: isoperimetry and processes* **23**. Springer Science & Business Media.
- [41] LI, S., XIE, Y., DAI, H. and SONG, L. (2015). M -Statistic for Kernel Change-Point Detection. *Advances in Neural Information Processing Systems* 3366–3374.
- [42] LIU, J., WU, S. and ZIDEK, J. V. (1997). On segmented multivariate regression. *Statistica Sinica* **7** 497–525.
- [43] LIU, S., SUZUKI, T., RELATOR, R., SESE, J., SUGIYAMA, M. and FUKUMIZU, K. (2017). Support consistency of direct sparse-change learning in Markov networks. *Annals of Statistics (accepted)*. arXiv:1407.0581.
- [44] PAGE, E. S. (1955). A test for a change in a parameter occurring at an unknown point. *Biometrika* 523–527.
- [45] RITOV, Y., RAZ, A. and BERGMAN, H. (2002). Detection of onset of neuronal activity by allowing for heterogeneity in the change points. *Journal of neuroscience methods* **122** 25–42.
- [46] SCHÖLKOPF, B. and SMOLA, A. J. (2002). *Learning with kernels: Support vector machines, regularization, optimization, and beyond*. MIT press.
- [47] SHAO, J. (1997). An asymptotic theory for linear model selection. *Statistica Sinica* **7** 221–242.
- [48] SHARIPOV, O., TEWES, J. and WENDLER, M. (2016). Sequential block bootstrap in a Hilbert space with application to change point analysis. *The Canadian Journal of Statistics. La Revue Canadienne de Statistique* **44** 300–322.
- [49] SPOKOINY, V. (2009). Multiscale local change point detection with applications to value-at-risk. *Annals of Statistics* 1405–1436.
- [50] SRIPERUMBUDUR, B. K., FUKUMIZU, K., GRETTON, A., LANCKRIET, G. R. G. and SCHÖLKOPF, B. (2009). Kernel Choice and Classifiability for RKHS Embeddings of Probability Distributions. In *Advances in Neural Information Processing Systems*, **21** NIPS Foundation.
- [51] TARTAKOVSKY, A., NIKIFOROV, I. V. and BASSEVILLE, M. (2014). *Sequential Analysis: Hypothesis Testing and Changepoint Detection*. *Monographs on Statistics and Applied Probability* **136**. Chapman and Hall/CRC, Boca Raton, FL.
- [52] VOGT, M. and DETTE, H. (2015). Detecting gradual changes in locally stationary processes. *Annals of Statistics* **43** 713–740.
- [53] WANG, T. and SAMWORTH, R. J. (2016). High-dimensional changepoint estimation via sparse projection. <https://arxiv.org/abs/1606.06246>.
- [54] YAO, Y.-C. (1988). Estimating the number of change-points via Schwarz’ criterion. *Statistics & Probability Letters* **6** 181–189.
- [55] YAO, Y.-C. and AU, S.-T. (1989). Least-squares estimation of a step function. *Sankhyā: The Indian Journal of Statistics, Series A* 370–381.
- [56] ZOU, C., YIN, G., FENG, L. and WANG, Z. (2014). Nonparametric maximum likelihood approach to multiple change-point problems. *Annals of Statistics* **42** 970–1002.